

(12) **United States Patent**  
**Archer et al.**

(10) **Patent No.:** **US 9,116,750 B2**  
(45) **Date of Patent:** **Aug. 25, 2015**

(54) **OPTIMIZING COLLECTIVE COMMUNICATIONS WITHIN A PARALLEL COMPUTER**

(75) Inventors: **Charles J. Archer**, Rochester, MN (US); **Michael A. Blocksome**, Rochester, MN (US); **Joseph D. Ratterman**, Seattle, WA (US); **Brian E. Smith**, Knoxville, TN (US)

(73) Assignee: **International Business Machines Corporation**, Armonk, NY (US)

(\*) Notice: Subject to any disclaimer, the term of this patent is extended or adjusted under 35 U.S.C. 154(b) by 398 days.

(21) Appl. No.: **13/569,614**

(22) Filed: **Aug. 8, 2012**

(65) **Prior Publication Data**

US 2014/0047451 A1 Feb. 13, 2014

(51) **Int. Cl.**  
**G06F 9/50** (2006.01)

(52) **U.S. Cl.**  
CPC ..... **G06F 9/5061** (2013.01); **G06F 2209/505** (2013.01)

(58) **Field of Classification Search**  
None  
See application file for complete search history.

(56) **References Cited**

U.S. PATENT DOCUMENTS

7,178,145	B2 *	2/2007	Bono	718/100
8,286,170	B2 *	10/2012	Franaszek	718/102
8,336,056	B1 *	12/2012	Gadir	718/104
2006/0242389	A1 *	10/2006	Browning et al.	712/229
2006/0253682	A1 *	11/2006	Armstrong et al.	711/173
2007/0150657	A1 *	6/2007	Yigzaw et al.	711/128
2008/0114973	A1 *	5/2008	Norton et al.	712/228

2008/0229321	A1 *	9/2008	Krieger et al.	718/104
2008/0235684	A1 *	9/2008	Anand et al.	718/100
2009/0031317	A1 *	1/2009	Gopalan et al.	718/103
2009/0031318	A1 *	1/2009	Gopalan et al.	718/103
2009/0187915	A1 *	7/2009	Chew et al.	718/104
2009/0320040	A1 *	12/2009	Robison	718/105
2010/0005080	A1 *	1/2010	Pike et al.	707/4
2010/0017804	A1 *	1/2010	Gupta et al.	718/102
2010/0115236	A1 *	5/2010	Bataineh et al.	712/29

(Continued)

OTHER PUBLICATIONS

Broquedis et al. "hwloc: a Generic Framework for managing Hardware Affinities in HPC Applications"; 2010 18<sup>th</sup> Euromicro Conference on Parallel, Distributed and Networked based Processing; (Broquedis\_2010.pdf; pp. 1-7).\*

(Continued)

*Primary Examiner* — Emerson Puentes

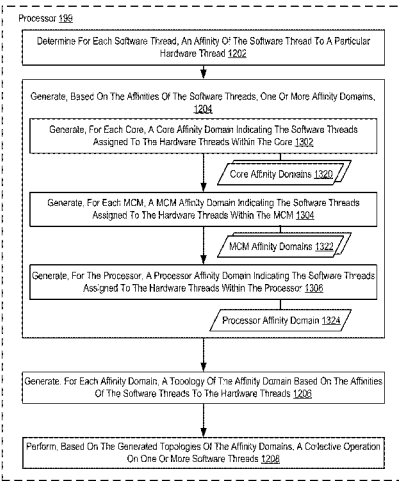
*Assistant Examiner* — Hiren Patel

(74) *Attorney, Agent, or Firm* — Edward J. Lenart; Kennedy Lenart Spraggins LLP

(57) **ABSTRACT**

Methods, apparatuses, and computer program products for optimizing collective communications within a parallel computer comprising a plurality of hardware threads for executing software threads of a parallel application are provided. Embodiments include a processor of a parallel computer determining for each software thread, an affinity of the software thread to a particular hardware thread. Each affinity indicates an assignment of a software thread to a particular hardware thread. The processor also generates one or more affinity domains based on the affinities of the software threads. Embodiments also include a processor generating, for each affinity domain, a topology of the affinity domain based on the affinities of the software threads to the hardware threads. According to embodiments of the present application, a processor also performs, based on the generated topologies of the affinity domains, a collective operation on one or more software threads.

**13 Claims, 17 Drawing Sheets**



(56)

**References Cited**

## U.S. PATENT DOCUMENTS

2010/0287279	A1 *	11/2010	Vaidyanathan et al. ....	709/226
2011/0035555	A1 *	2/2011	Wang et al. ....	711/137
2011/0088038	A1 *	4/2011	Kruglick ....	718/104
2011/0161969	A1 *	6/2011	Arndt et al. ....	718/103
2012/0054750	A1 *	3/2012	Saripalli ....	718/100
2012/0102500	A1 *	4/2012	Waddington et al. ....	718/104
2012/0102501	A1 *	4/2012	Waddington et al. ....	718/105
2012/0272016	A1 *	10/2012	Bello et al. ....	711/154
2014/0013330	A1 *	1/2014	Wang et al. ....	718/103
2014/0143783	A1 *	5/2014	Bose et al. ....	718/102

## OTHER PUBLICATIONS

Brochard et al., "IBM Tuning Guide for High Performance Computing Applications"; IBM Systems and Technology Group; Release 1.0; Apr. 7, 2009 (Brochard\_APR2009.pdf; pp. 1-41).\*

Engelberts et al.; "Best Practice Guide—IBM Power"; Jun. 19, 2012; (Engelberts\_IBM\_JUN2012.pdf; pp. 1-52).\*

"A Performance Guide for HPC Applications on the IBM System x iDataPlex dx360 M4 System"; IBM, Release 1.0.2, Jun. 19, 2012 (IBM\_JUN2012.pdf; pp. 1-153).\*

Srinivasan et al., "Processor Affinity and MPI Performance on SMP-CMP Clusters", The 11th IPDPS Workshop on Parallel and Distributed Scientific and Engineering Computing ('PDSEC'), Apr. 2010, pp. 1-8, IEEE.org, USA.

IBM, "IBM Parallel Environment Runtime Edition for AIX, V1.1 enables parallel application development and execution on select IBM Power Systems clusters", www.ibm.com [online] Jul. 2011, pp. 1-14, [accessed online Feb. 21, 2012], URL: <http://www-01.ibm.com/common/ssi/cgi-bin/ssialias?subtype=ca&infotype=an&appname=iSource&supplier=897&letternum=ENUS211-237>.

\* cited by examiner

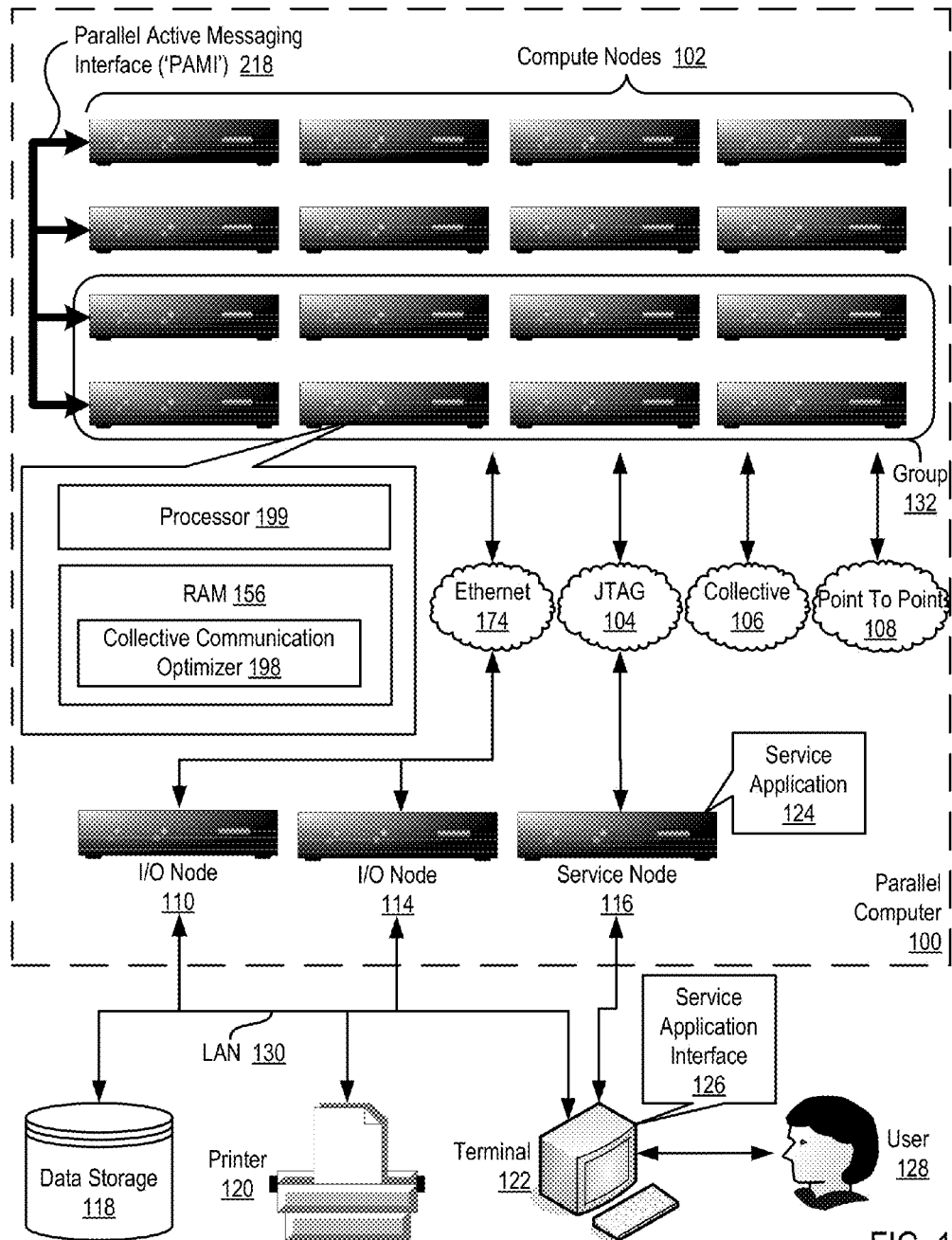


FIG. 1

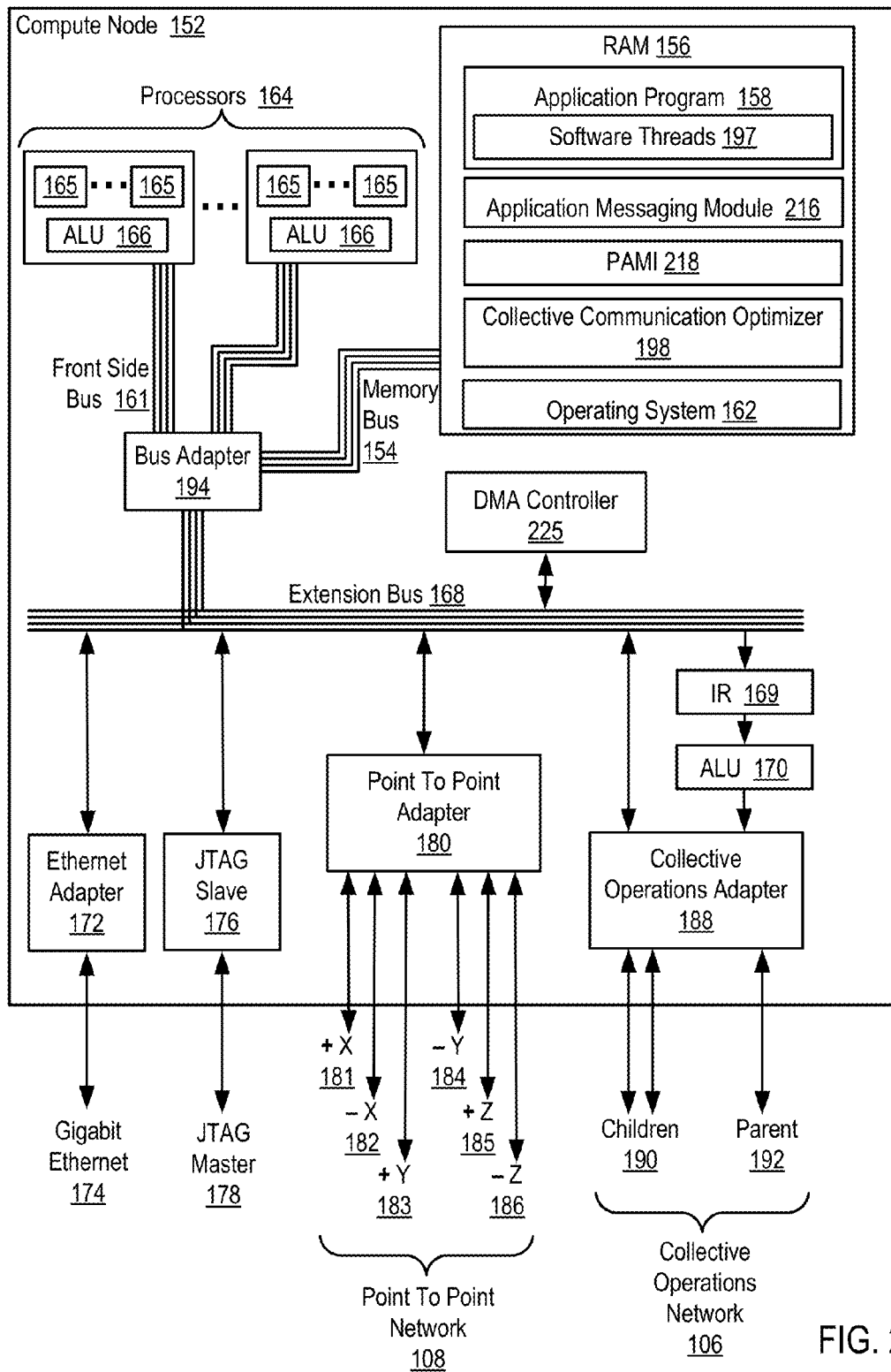


FIG. 2

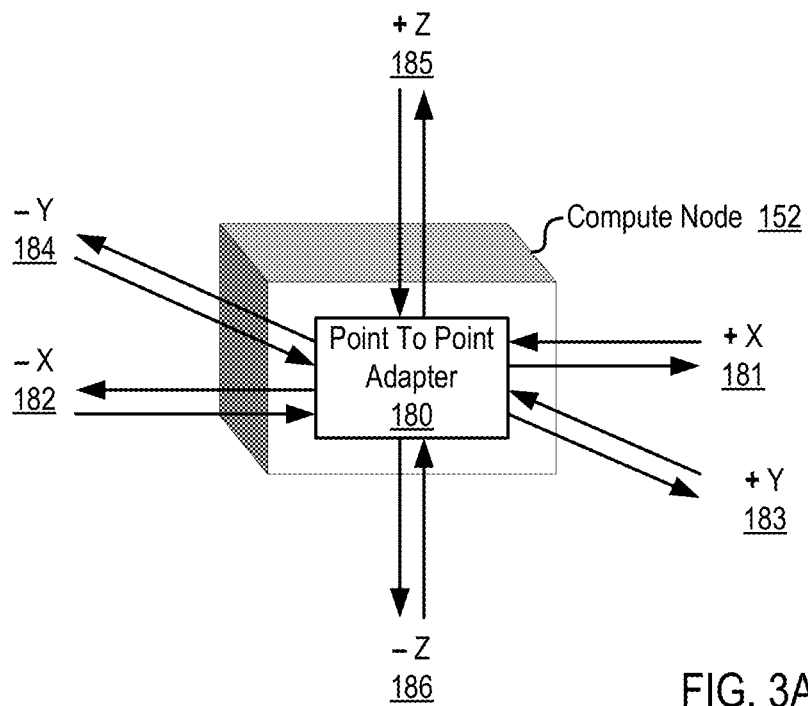
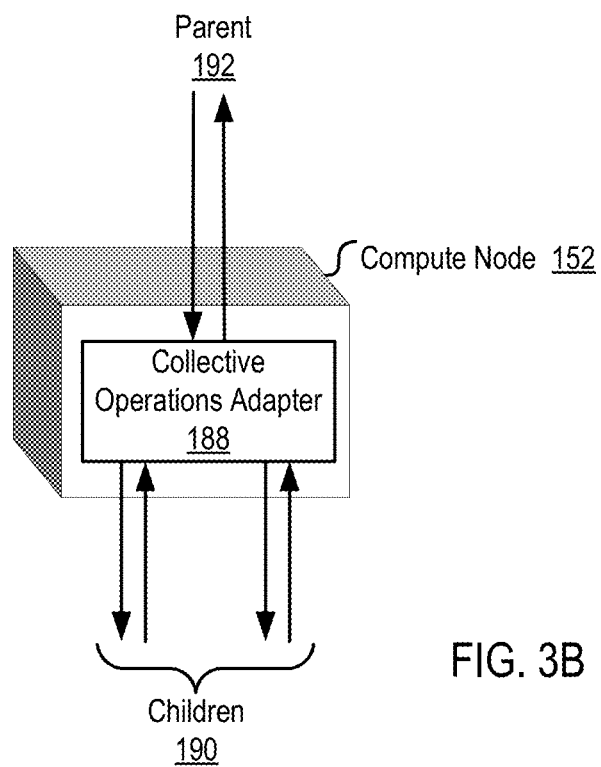
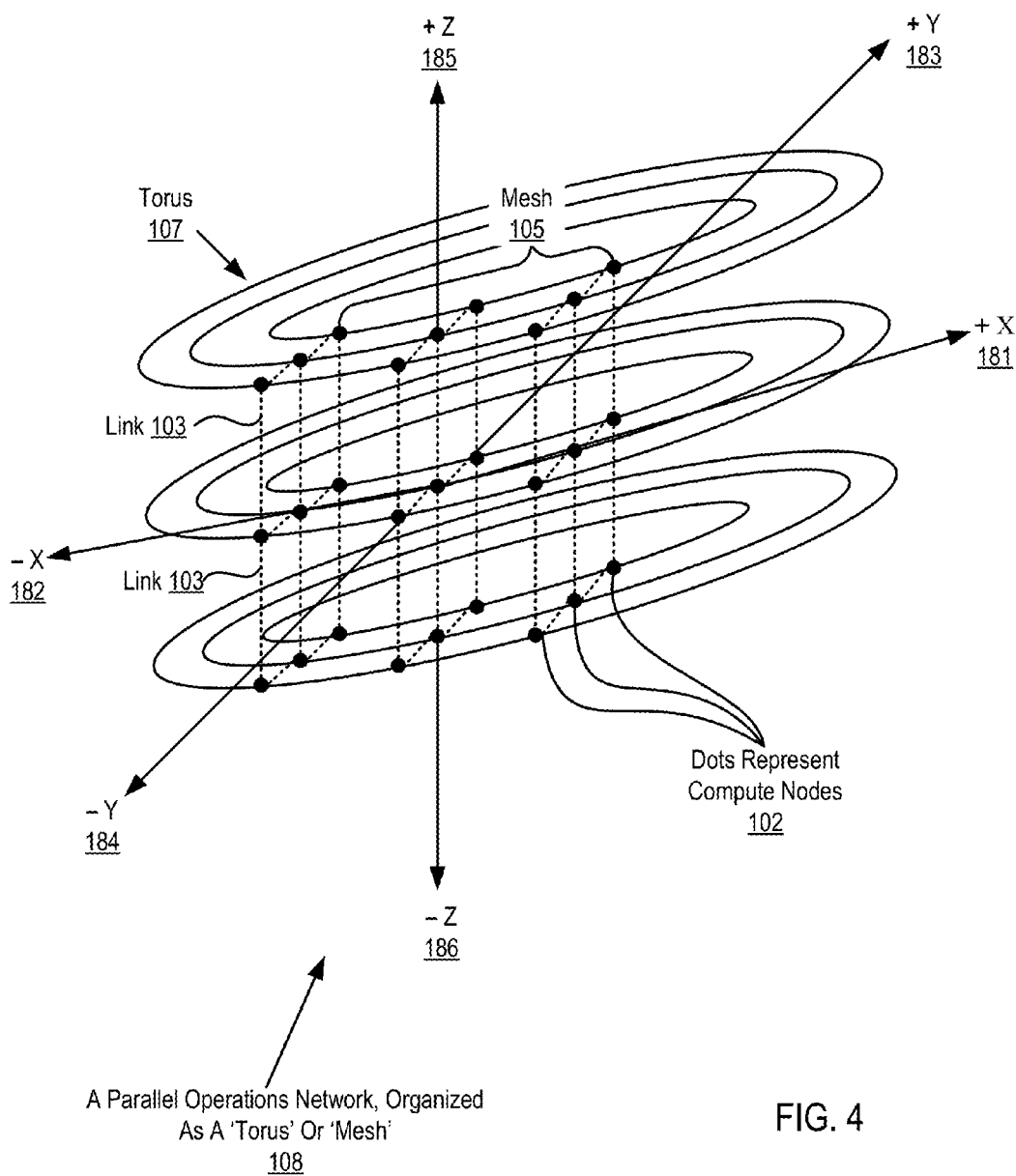


FIG. 3A





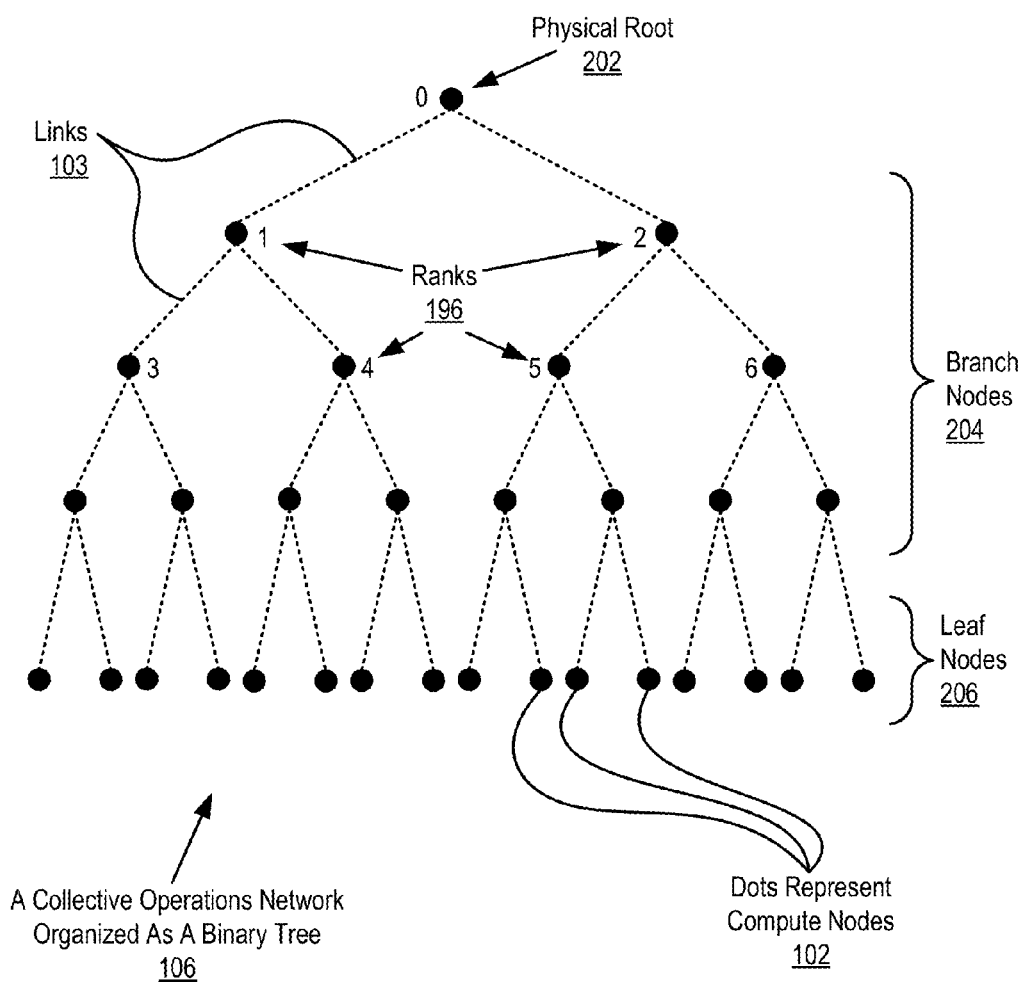


FIG. 5



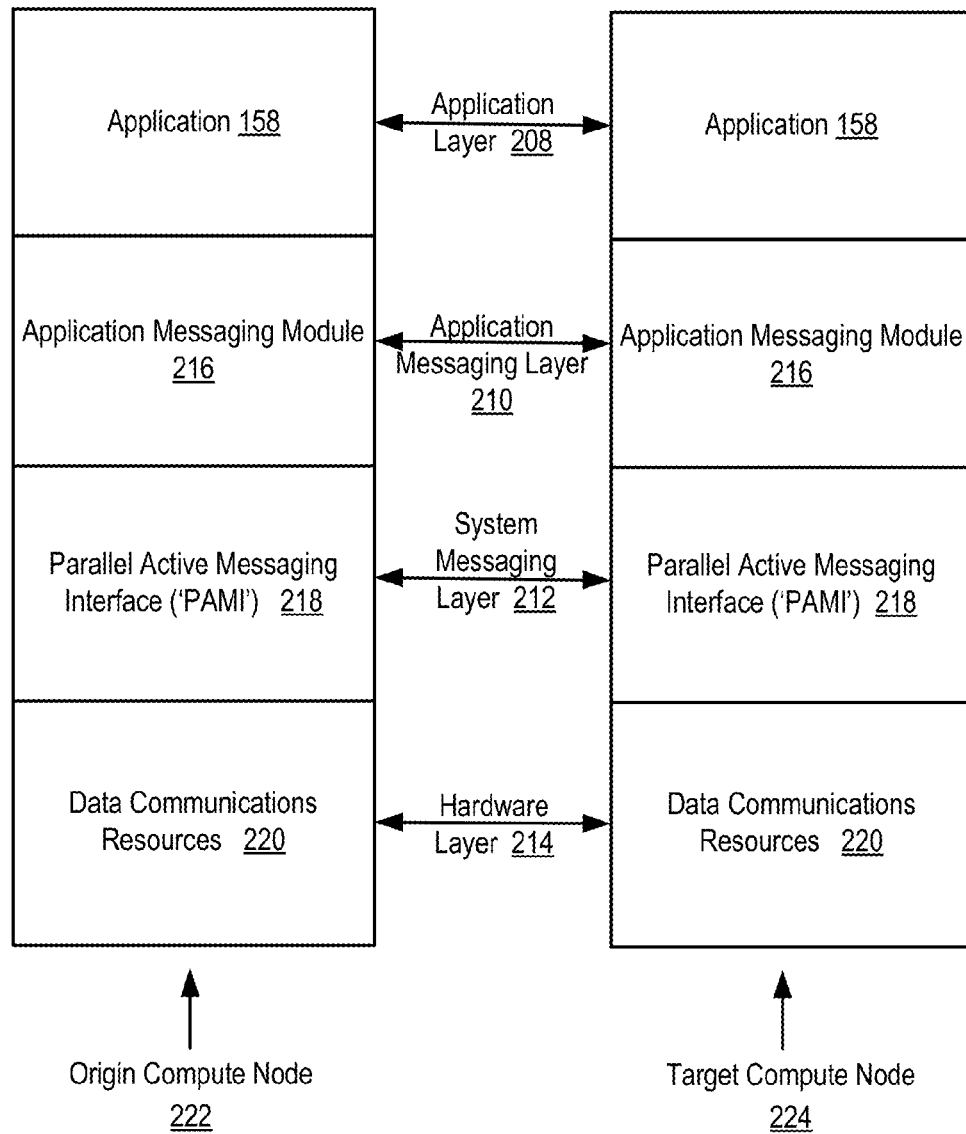


FIG. 6

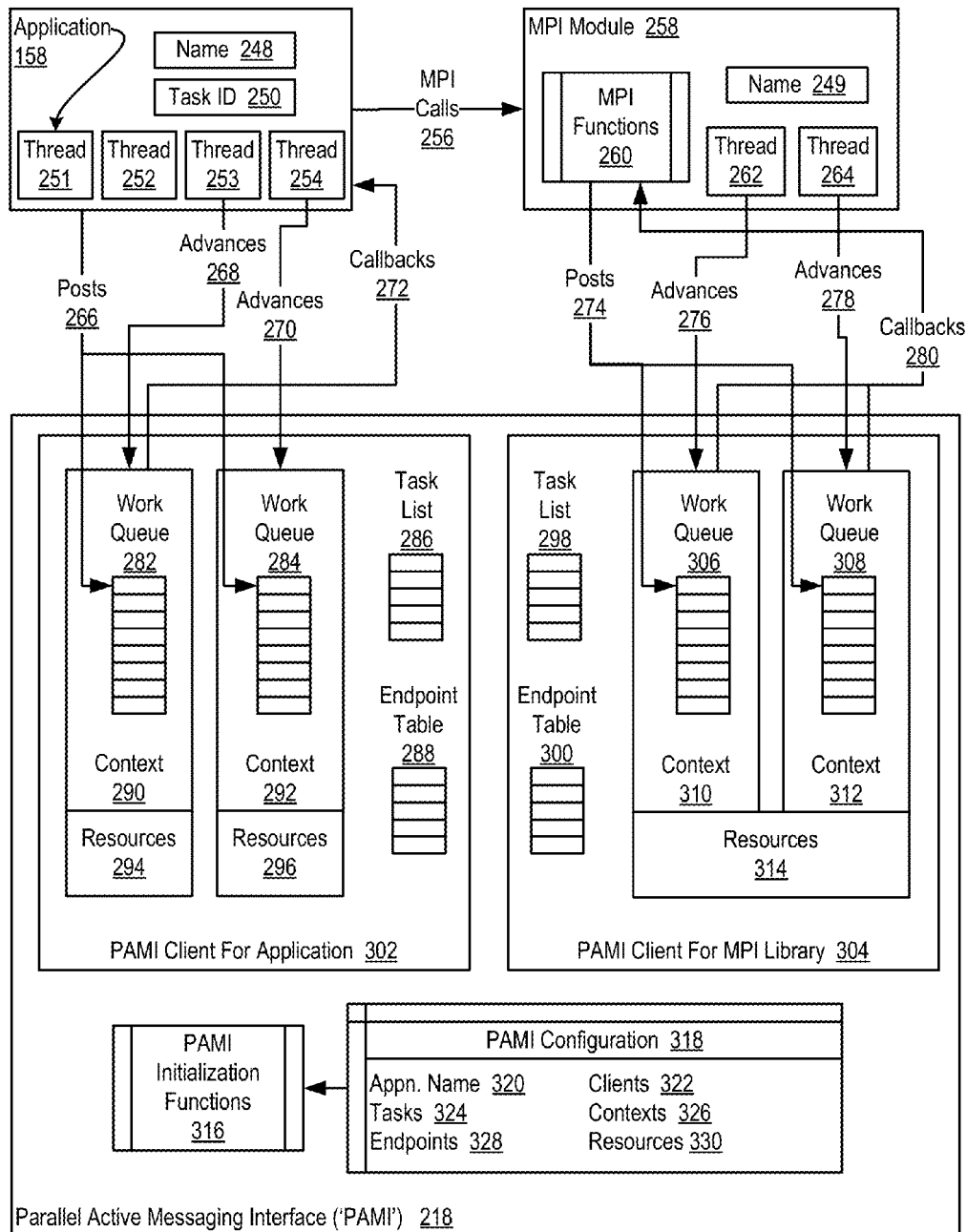
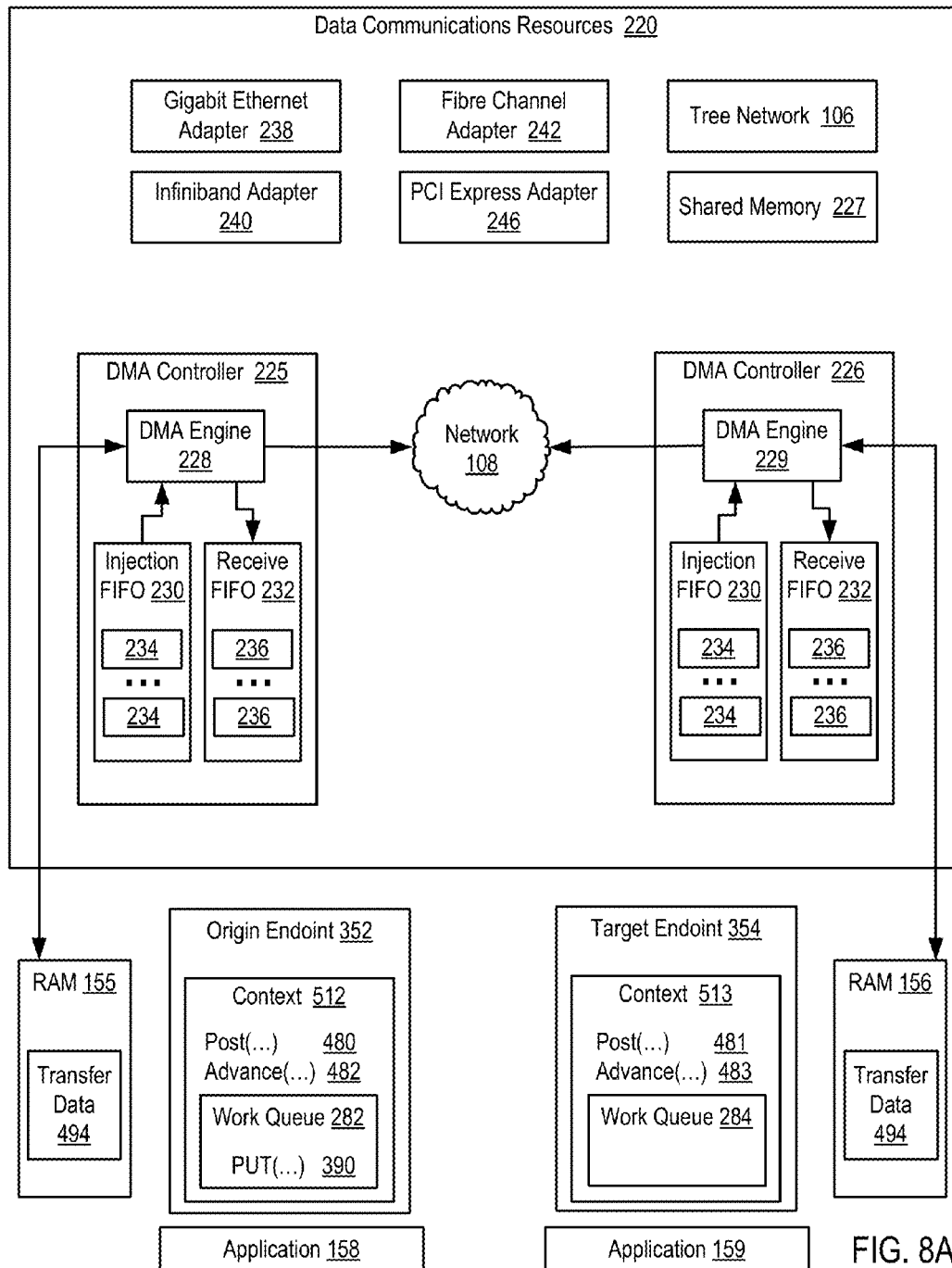


FIG. 7



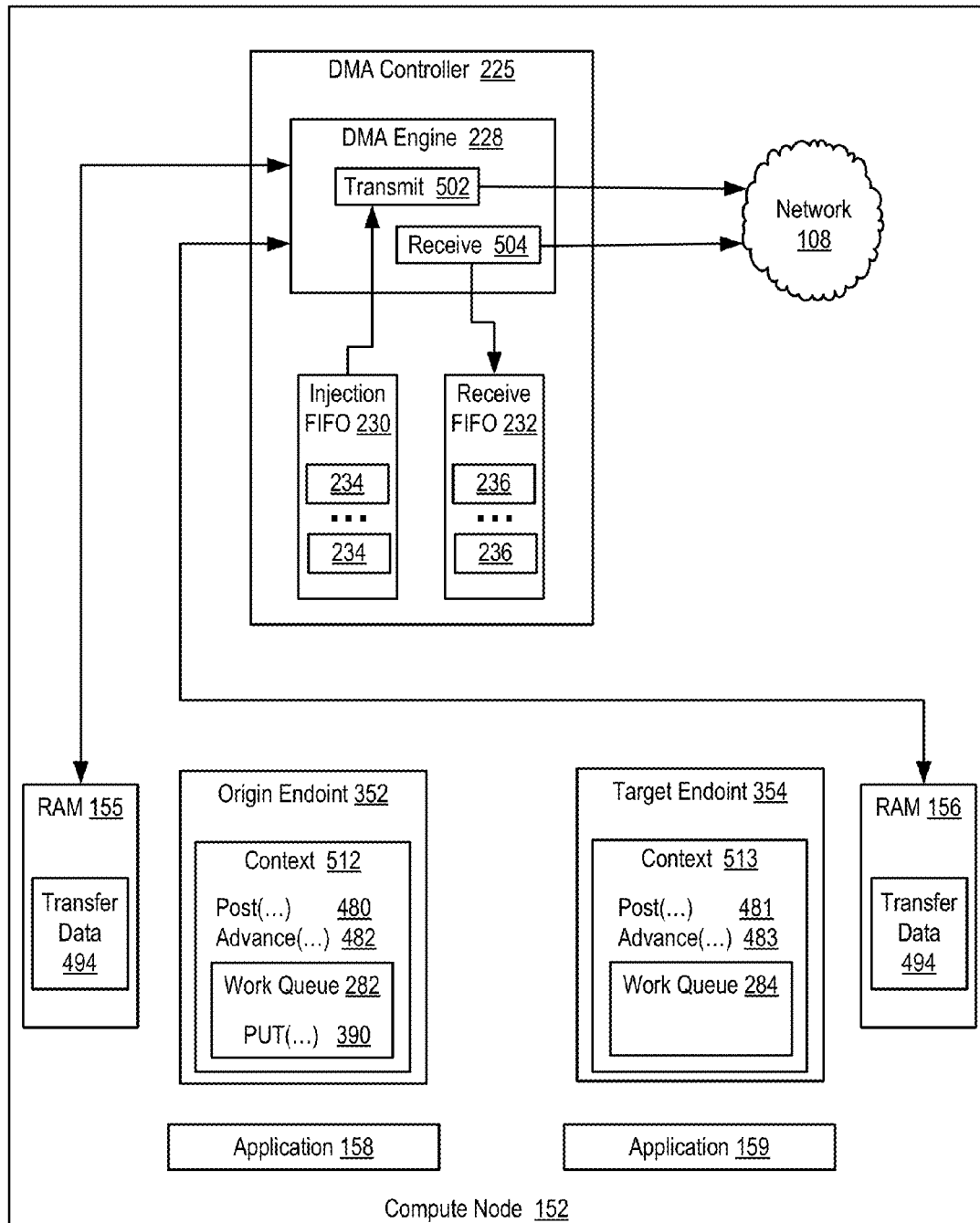


FIG. 8B

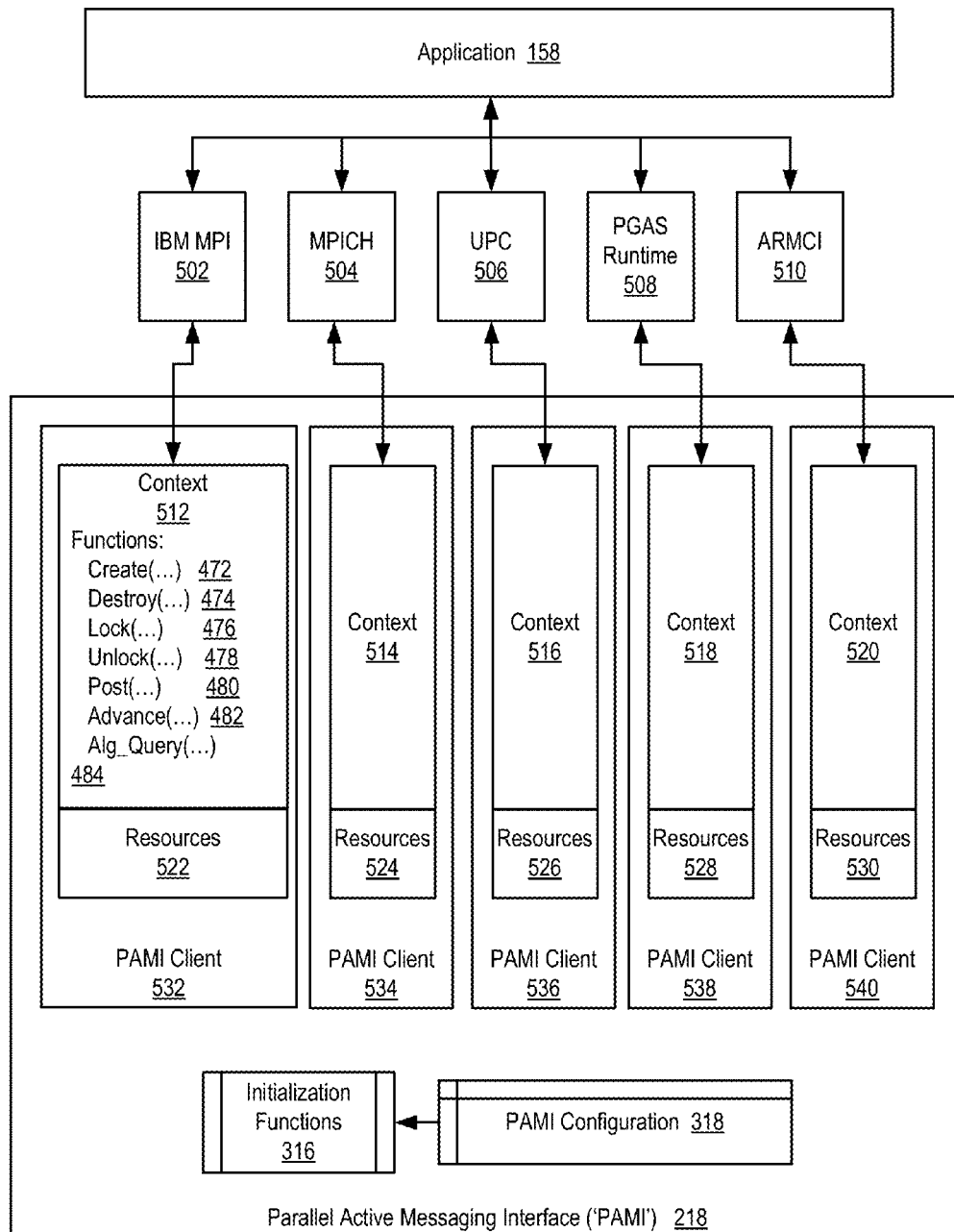


FIG. 9

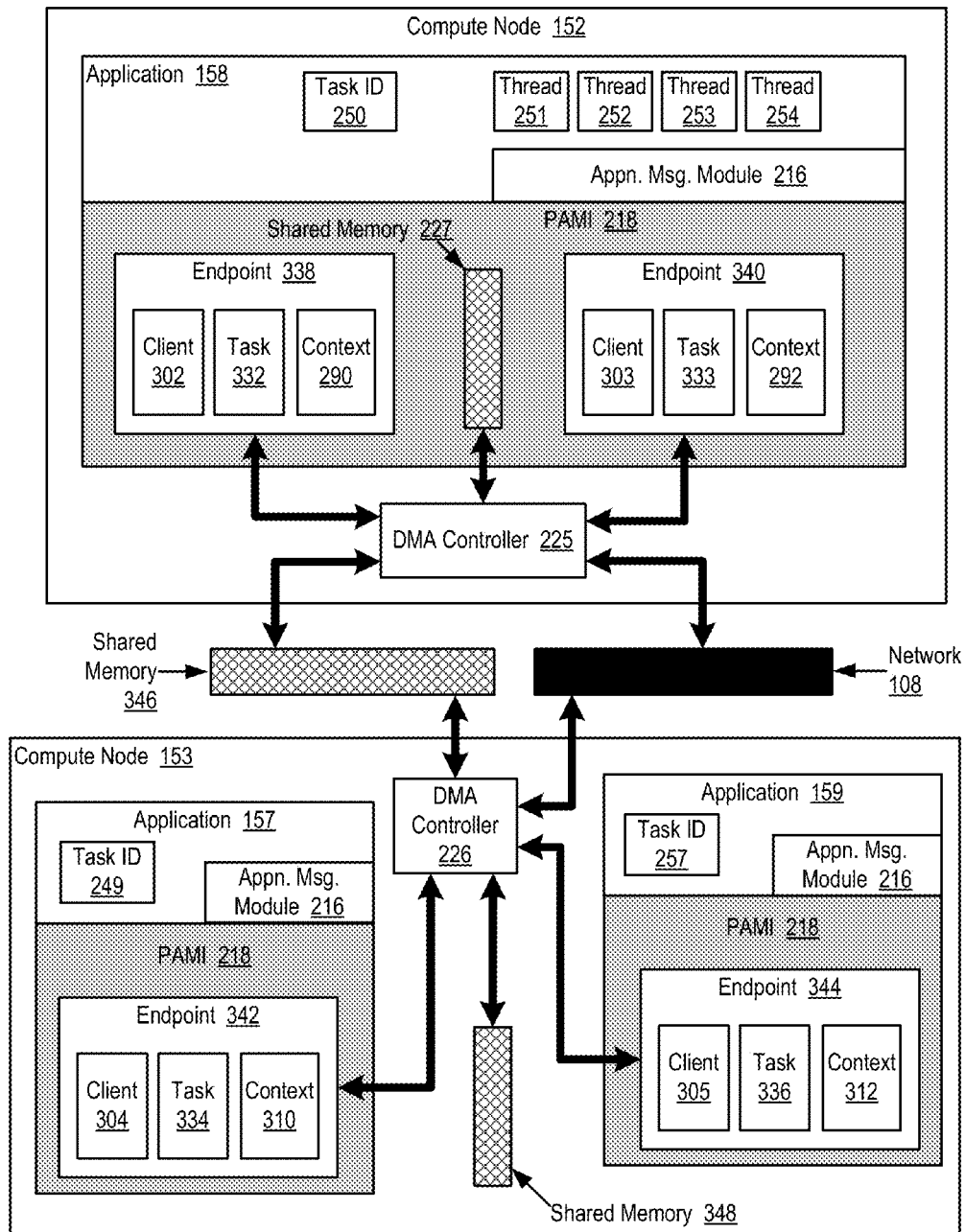


FIG. 10

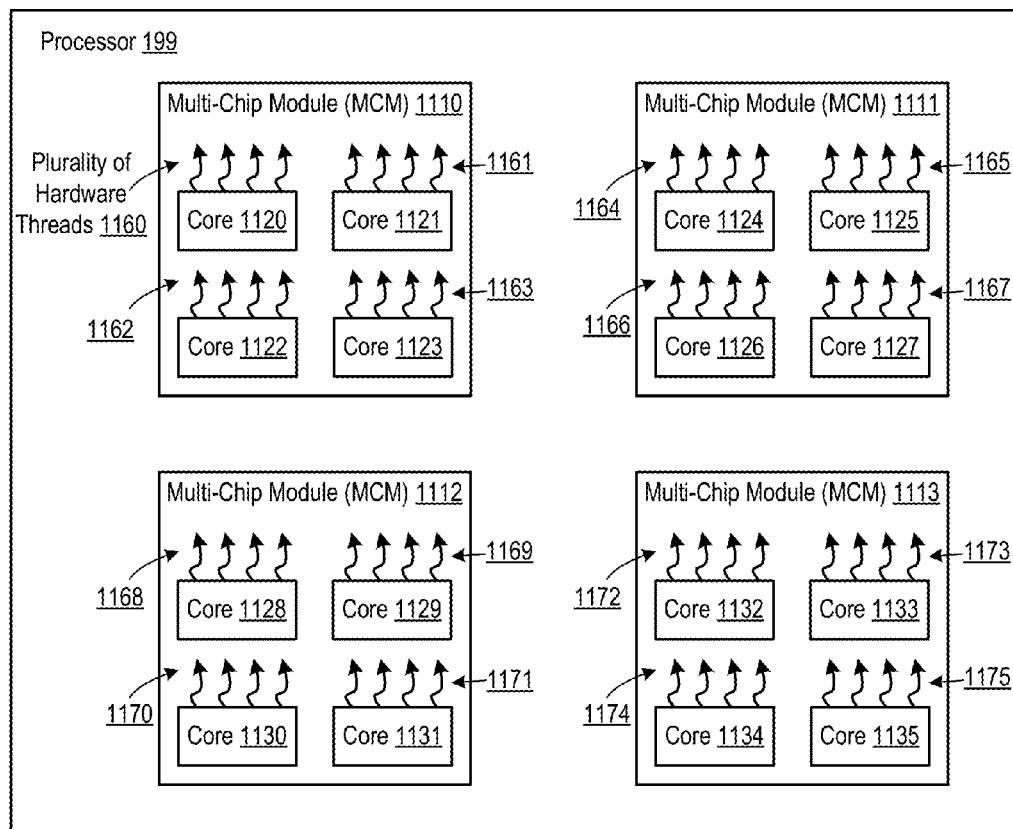


FIG. 11

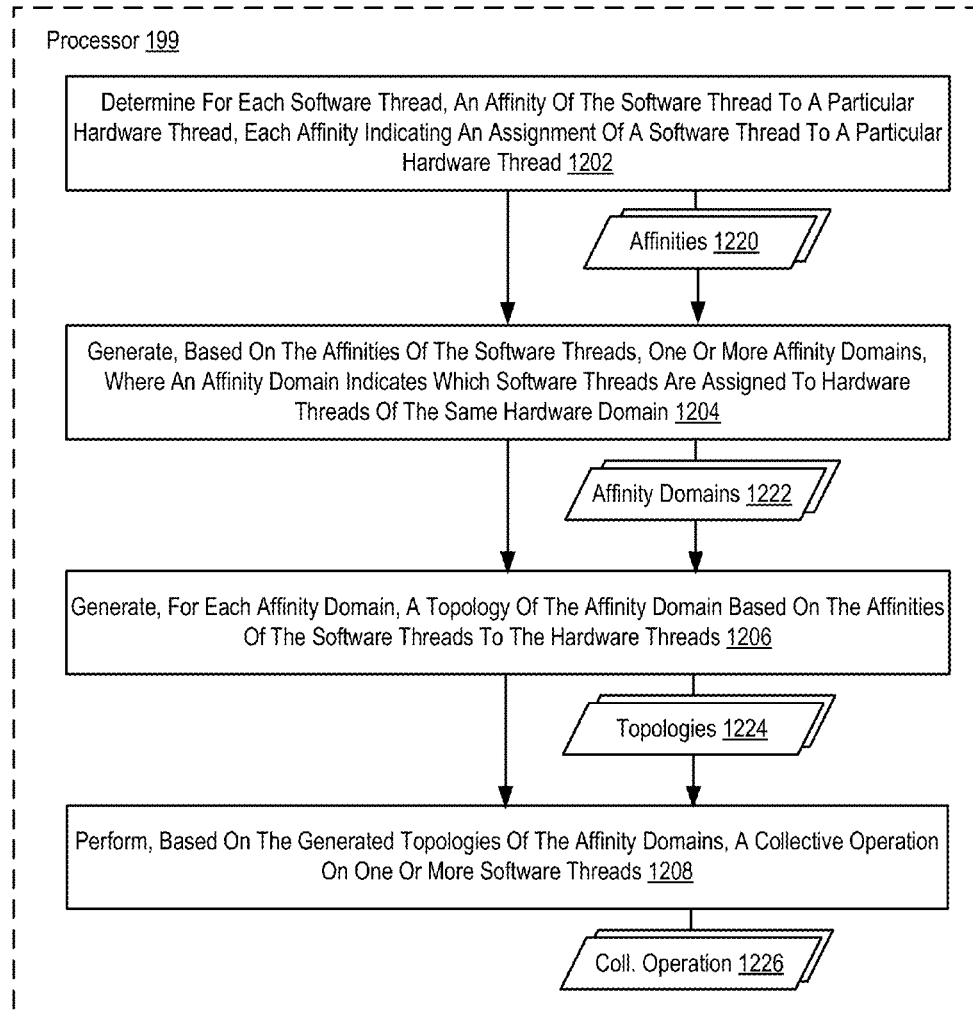


FIG. 12



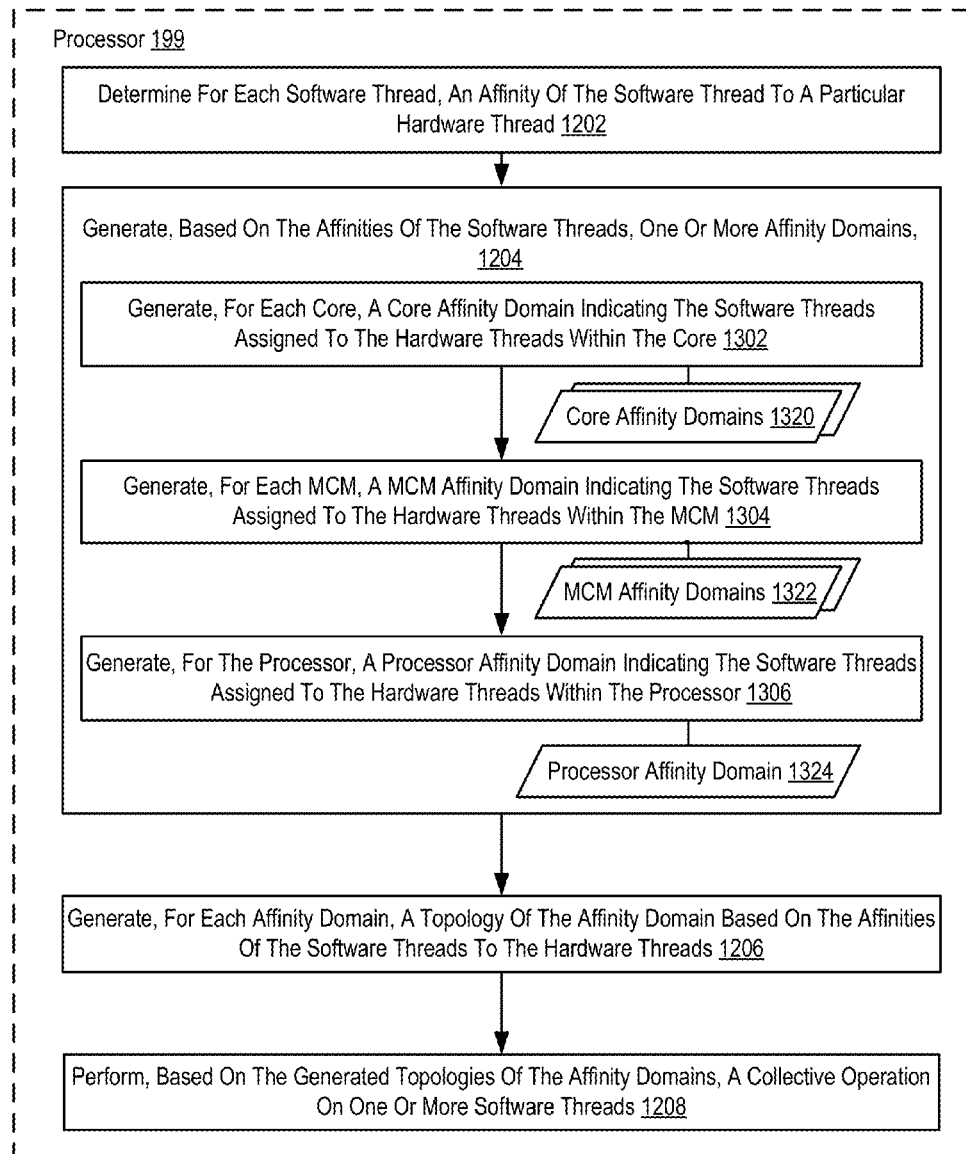


FIG. 13

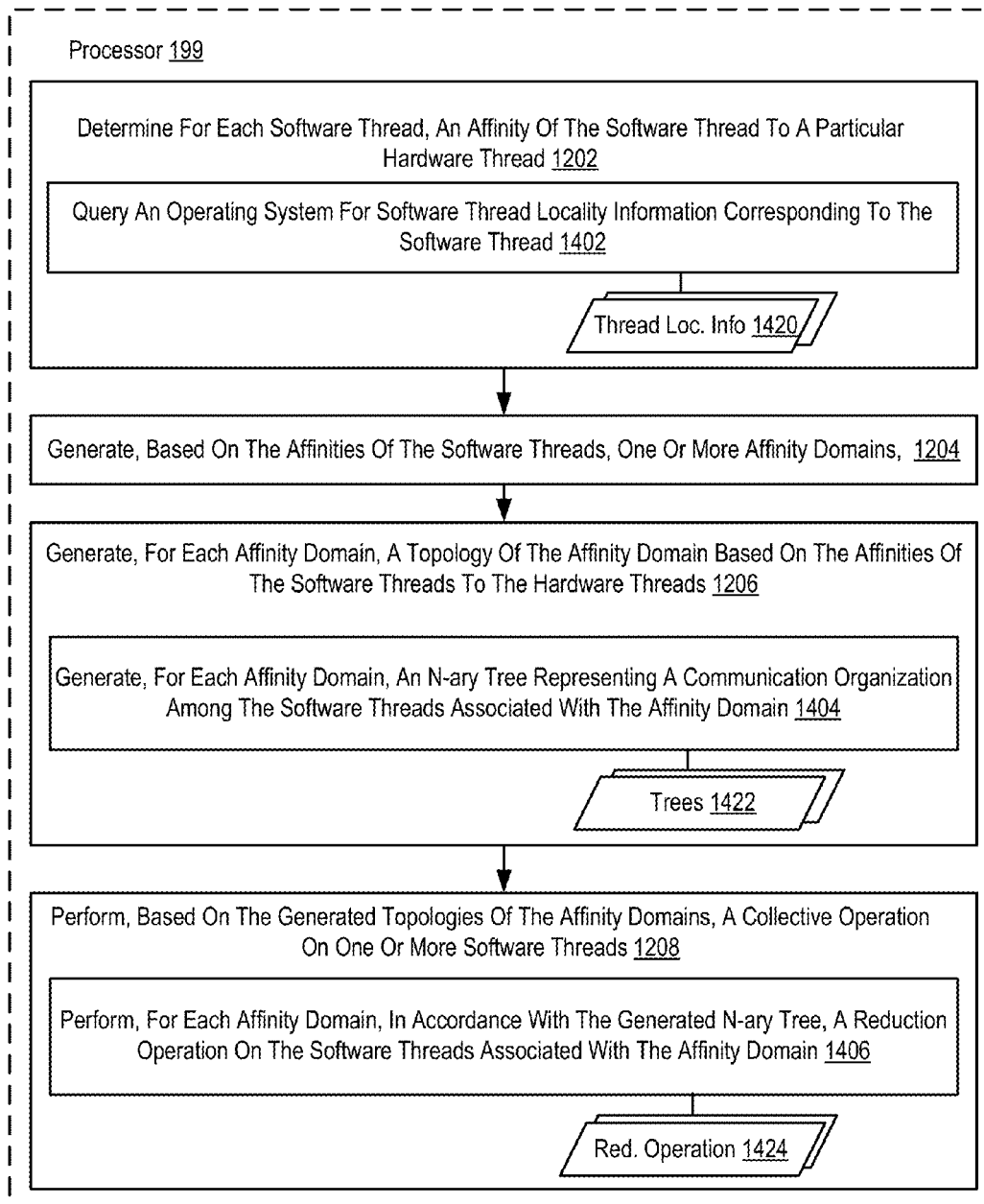


FIG. 14

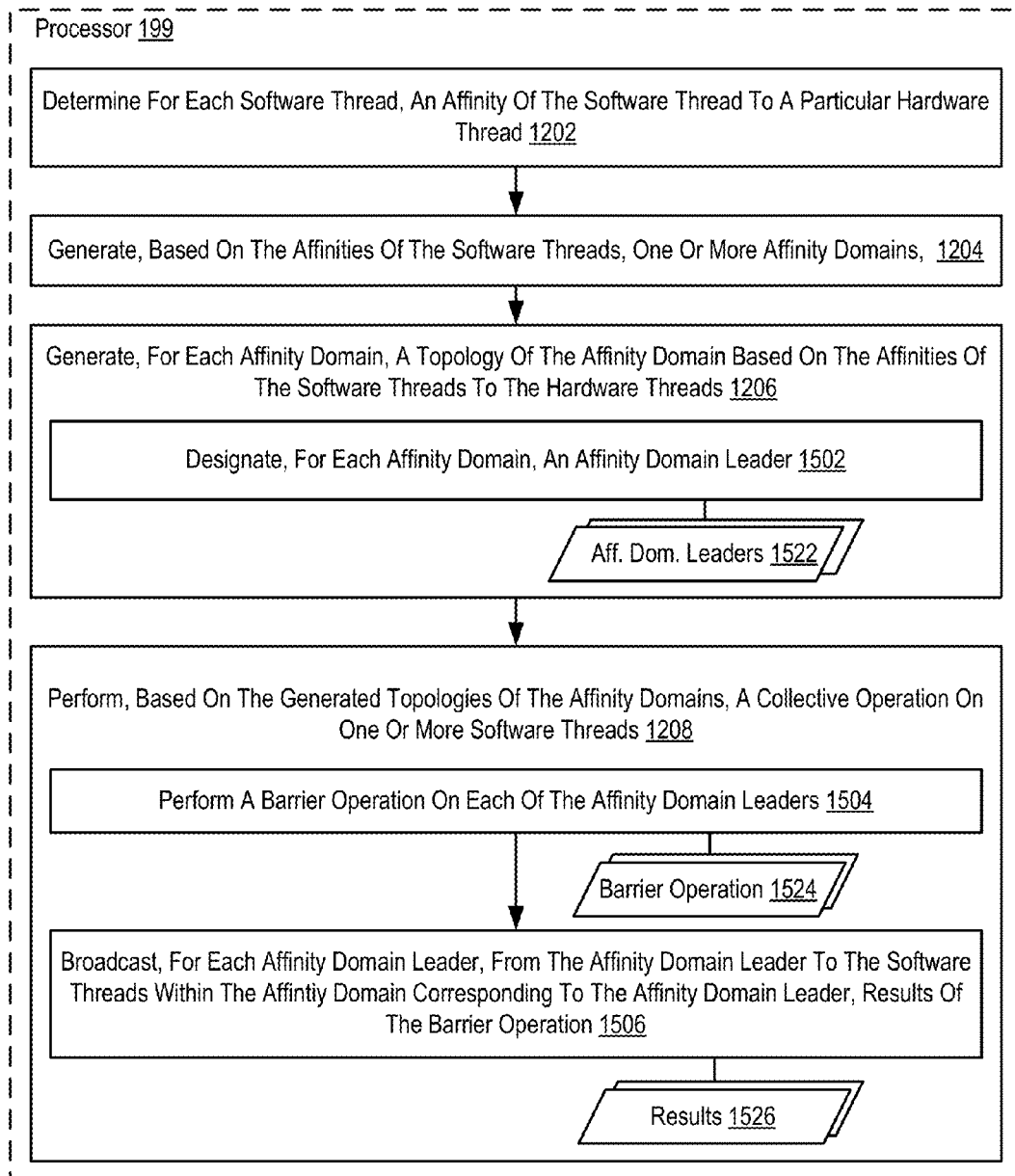


FIG. 15

# OPTIMIZING COLLECTIVE COMMUNICATIONS WITHIN A PARALLEL COMPUTER

## BACKGROUND OF THE INVENTION

### 1. Field of the Invention

The field of the invention is data processing, or, more specifically, methods, apparatuses, and computer program products for optimizing collective communications within a parallel computer.

### 2. Description of Related Art

The development of the EDVAC computer system of 1948 is often cited as the beginning of the computer era. Since that time, computer systems have evolved into extremely complicated devices. Today's computers are much more sophisticated than early systems such as the EDVAC. Computer systems typically include a combination of hardware and software components, application programs, operating systems, processors, buses, memory, input/output devices, and so on. As advances in semiconductor processing and computer architecture push the performance of the computer higher and higher, more sophisticated computer software has evolved to take advantage of the higher performance of the hardware, resulting in computer systems today that are much more powerful than just a few years ago.

Parallel computing is an area of computer technology that has experienced advances. Parallel computing is the simultaneous execution of the same application (split up and specially adapted) on multiple processors in order to obtain results faster. Parallel computing is based on the fact that the process of solving a problem usually can be divided into smaller jobs, which may be carried out simultaneously with some coordination.

Parallel computers execute parallel algorithms. A parallel algorithm can be split up to be executed a piece at a time on many different processing devices, and then put back together again at the end to get a data processing result. Some algorithms are easy to divide up into pieces. Splitting up the job of checking all of the numbers from one to a hundred thousand to see which are primes could be done, for example, by assigning a subset of the numbers to each available processor, and then putting the list of positive results back together. In this specification, the multiple processing devices that execute the individual pieces of a parallel program are referred to as 'compute nodes.' A parallel computer is composed of compute nodes and other processing nodes as well, including, for example, input/output ('I/O') nodes, and service nodes.

Parallel algorithms are valuable because it is faster to perform some kinds of large computing jobs via a parallel algorithm than it is via a serial (non-parallel) algorithm, because of the way modern processors work. It is far more difficult to construct a computer with a single fast processor than one with many slow processors with the same throughput. There are also certain theoretical limits to the potential speed of serial processors. On the other hand, every parallel algorithm has a serial part and so parallel algorithms have a saturation point. After that point adding more processors does not yield any more throughput but only increases the overhead and cost.

Parallel algorithms are designed also to optimize one more resource the data communications requirements among the nodes of a parallel computer. There are two ways parallel processors communicate, shared memory or message passing. Shared memory processing needs additional locking for the data and imposes the overhead of additional processor and

bus cycles and also serializes some portion of the algorithm. Message passing processing uses high-speed data communications networks and message buffers, but this communication adds transfer overhead on the data communications networks as well as additional memory need for message buffers and latency in the data communications among nodes. Designs of parallel computers use specially designed data communications links so that the communication overhead will be small but it is the parallel algorithm that decides the volume of the traffic.

Many data communications network architectures are used for message passing among nodes in parallel computers. Compute nodes may be organized in a network as a 'torus' or 'mesh,' for example. Also, compute nodes may be organized in a network as a tree. A torus network connects the nodes in a three-dimensional mesh with wrap around links. Every node is connected to its six neighbors through this torus network, and each node is addressed by its x, y, z coordinate in the mesh. In a tree network, the nodes typically are connected into a binary tree: each node has a parent and two children (although some nodes may only have zero children or one child, depending on the hardware configuration). In computers that use a torus and a tree network, the two networks typically are implemented independently of one another, with separate routing circuits, separate physical links, and separate message buffers.

A torus network lends itself to point to point operations, but a tree network typically is inefficient in point to point communication. A tree network, however, does provide high bandwidth and low latency for certain collective operations, message passing operations where all compute nodes participate simultaneously, such as, for example, an allgather.

## SUMMARY OF THE INVENTION

Methods, apparatuses, and computer program products for optimizing collective communications within a parallel computer comprising a plurality of hardware threads for executing software threads of a parallel application are provided. Embodiments include a processor of a parallel computer determining for each software thread, an affinity of the software thread to a particular hardware thread. Each affinity indicates an assignment of a software thread to a particular hardware thread. The processor also generates one or more affinity domains based on the affinities of the software threads. An affinity domain indicates which software threads are assigned to hardware threads of the same hardware domain. A hardware domain may refer to an organizational level within a processor. Examples of a hardware domain include a core domain, a multi-chip module (MCM) domain, and a processor domain. Embodiments also include a processor generating, for each affinity domain, a topology of the affinity domain based on the affinities of the software threads to the hardware threads. According to embodiments of the present application, a processor also performs, based on the generated topologies of the affinity domains, a collective operation on one or more software threads.

The foregoing and other objects, features and advantages of the invention will be apparent from the following more particular descriptions of exemplary embodiments of the invention as illustrated in the accompanying drawings wherein like reference numbers generally represent like parts of exemplary embodiments of the invention.

## BRIEF DESCRIPTION OF THE DRAWINGS

FIG. 1 sets forth a block and network diagram of an example parallel computer that optimizes collective communications according to embodiments of the present invention.

FIG. 2 sets forth a block diagram of an example compute node for use in a parallel computer that optimizes collective communications according to embodiments of the present invention.

FIG. 3A illustrates an example of a Point To Point Adapter useful in parallel computers that optimize collective communications according to embodiments of the present invention.

FIG. 3B illustrates an example of a Collective Operations Adapter useful in a parallel computer that optimizes collective communications according to embodiments of the present invention.

FIG. 4 sets forth a line drawing illustrating an example data communications network optimized for point-to-point operations useful in parallel computers that optimize collective communications according to embodiments of the present invention.

FIG. 5 illustrates an example data communications network optimized for collective operations by organizing compute nodes in a tree.

FIG. 6 sets forth a block diagram of an example protocol stack useful in parallel computers that optimize collective communications according to embodiments of the present invention.

FIG. 7 sets forth a functional block diagram of an example PAMI for use in parallel computers that optimize collective communications according to embodiments of the present invention.

FIG. 8A sets forth a block diagram of example data communications resources useful in parallel computers that optimize collective communications according to embodiments of the present invention.

FIG. 8B sets forth a functional block diagram of an example DMA controller operatively coupled to a network—in an architecture where this DMA controller is the only DMA controller on a compute node—and an origin endpoint and its target endpoint are both located on the same compute node.

FIG. 9 sets forth a functional block diagram of an example PAMI useful in parallel computers that optimize collective communications according to embodiments of the present invention in which the example PAMI operates, on behalf of an application, with multiple application messaging modules simultaneously.

FIG. 10 sets forth a functional block diagram of example endpoints useful in parallel computers that optimize collective communications according to embodiments of the present invention.

FIG. 11 sets forth a flow chart illustrating a functional block diagram of an example processor that includes collective communication optimizer useful for optimization of collective communications in a parallel computer according to embodiments of the present invention.

FIG. 12 sets forth a flow chart illustrating an example method of optimizing collective communications within a parallel computer according to embodiments of the present invention.

FIG. 13 sets forth a flow chart illustrating another example method of optimizing collective communications within a parallel computer according to embodiments of the present invention.

FIG. 14 sets forth a flow chart illustrating another example method of optimizing collective communications within a parallel computer according to embodiments of the present invention.

FIG. 15 sets forth a flow chart illustrating another example method of optimizing collective communications within a parallel computer according to embodiments of the present invention.

#### DETAILED DESCRIPTION OF EXEMPLARY EMBODIMENTS

Example methods, computers, and computer program products for optimizing collective communications within a parallel computer according to embodiments of the present invention are described with reference to the accompanying drawings, beginning with FIG. 1. FIG. 1 sets forth a block and network diagram of an example parallel computer (100) that optimizes collective communications according to embodiments of the present invention. The parallel computer (100) in the example of FIG. 1 is coupled to non-volatile memory for the computer in the form of data storage device (118), an output device for the computer in the form of printer (120), and an input/output device for the computer in the form of computer terminal (122).

The parallel computer (100) in the example of FIG. 1 includes a plurality of compute nodes (102). The compute nodes (102) are coupled for data communications by several independent data communications networks including a high speed Ethernet network (174), a Joint Test Action Group ('JTAG') network (104), a tree network (106) which is optimized for collective operations, and a torus network (108) which is optimized point to point operations. Tree network (106) is a data communications network that includes data communications links connected to the compute nodes so as to organize the compute nodes as a tree. Each data communications network is implemented with data communications links among the compute nodes (102). The data communications links provide data communications for parallel operations among the compute nodes of the parallel computer.

In addition, the compute nodes (102) of parallel computer (100) are organized into at least one operational group (132) of compute nodes for collective parallel operations on parallel computer (100). An operational group of compute nodes is the set of compute nodes upon which a collective parallel operation executes. Collective operations are implemented with data communications among the compute nodes of an operational group. Collective operations are those functions that involve all the compute nodes of an operational group. A collective operation is an operation, a message-passing computer program instruction that is executed simultaneously, that is, at approximately the same time, by all the compute nodes in an operational group of compute nodes. Such an operational group may include all the compute nodes in a parallel computer (100) or a subset all the compute nodes. Collective operations are often built around point to point operations. A collective operation requires that all processes on all compute nodes within an operational group call the same collective operation with matching arguments. A 'broadcast' is an example of a collective operation for moving data among compute nodes of an operational group. A 'reduce' operation is an example of a collective operation that executes arithmetic or logical functions on data distributed among the compute nodes of an operational group. A 'barrier operation' is an example of a collective operation that is issued by each thread participating in the barrier operation and only can complete when all threads have entered and left the collective barrier operation. An operational group may be implemented as, for example, an MPI 'communicator.'

'MPI' refers to 'Message Passing Interface,' a prior art applications messaging module or parallel communications

library, an application-level messaging module of computer program instructions for data communications on parallel computers. Such an application messaging module is disposed in an application messaging layer in a data communications protocol stack. Examples of prior-art parallel communications libraries that may be improved for use with parallel computers that optimize collective communications according to embodiments of the present invention include IBM's MPI library, the 'Parallel Virtual Machine' ('PVM') library, MPICH, OpenMPI, and LAM/MPI. MPI is promulgated by the MPI Forum, an open group with representatives from many organizations that define and maintain the MPI standard. MPI at the time of this writing is a de facto standard for communication among compute nodes running a parallel program on a distributed memory parallel computer. This specification sometimes uses MPI terminology for ease of explanation, although the use of MPI as such is not a requirement or limitation of the present invention.

Most collective operations are variations or combinations of four basic operations: broadcast, gather, scatter, and reduce. In a broadcast operation, all processes specify the same root process, whose buffer contents will be sent. Processes other than the root specify receive buffers. After the operation, all buffers contain the message from the root process.

A scatter operation, like the broadcast operation, is also a one-to-many collective operation. All processes specify the same receive count. The send arguments are only significant to the root process, whose buffer actually contains sendcount\*N elements of a given datatype, where N is the number of processes in the given group of compute nodes. The send buffer will be divided equally and dispersed from the root to all processes (including the root). Each process is assigned a sequential identifier termed a 'rank.' After the operation, the root has sent sendcount data elements to each process in increasing rank order. Rank 0 (the root process) receives the first sendcount data elements from the send buffer. Rank 1 receives the second sendcount data elements from the send buffer, and so on.

A gather operation is a many-to-one collective operation that is a complete reverse of the description of the scatter operation. That is, a gather is a many-to-one collective operation in which elements of a datatype are gathered from the ranked processes into a receive buffer of the root process.

A reduce operation is also a many-to-one collective operation that includes an arithmetic or logical function performed on two data elements. All processes specify the same 'count' and the same arithmetic or logical function. After the reduction, all processes have sent count data elements from compute node send buffers to the root process. In a reduction operation, data elements from corresponding send buffer locations are combined pair-wise by arithmetic or logical operations to yield a single corresponding element in the root process's receive buffer. Application specific reduction operations can be defined at runtime. Parallel communications libraries may support predefined operations. MPI, for example, provides the following pre-defined reduction operations:

MPI_MAX	maximum
MPI_MIN	minimum
MPI_SUM	sum
MPI_PROD	product
MPI_LAND	logical AND
MPI_BAND	bitwise AND
MPI_LOR	logical OR

-continued

MPI_BOR	bitwise OR
MPI_LXOR	logical exclusive OR
MPI_BXOR	bitwise exclusive OR

In addition to compute nodes, the example parallel computer (100) includes input/output ('I/O') nodes (110, 114) coupled to compute nodes (102) through one of the data communications networks (174). The I/O nodes (110, 114) provide I/O services between compute nodes (102) and I/O devices (118, 120, 122). I/O nodes (110, 114) are connected for data communications I/O devices (118, 120, 122) through local area network ('LAN') (130). Computer (100) also includes a service node (116) coupled to the compute nodes through one of the networks (104). Service node (116) provides service common to pluralities of compute nodes, loading programs into the compute nodes, starting program execution on the compute nodes, retrieving results of program operations on the compute nodes, and so on. Service node (116) runs a service application (124) and communicates with users (128) through a service application interface (126) that runs on computer terminal (122). As the term is used here, a parallel active messaging interface or 'PAMI' (218) is a system-level messaging layer in a protocol stack of a parallel computer that is composed of data communications endpoints each of which is specified with data communications parameters for a thread of execution on a compute node of the parallel computer. The PAMI is a 'parallel' interface in that many instances of the PAMI operate in parallel on the compute nodes of a parallel computer. The PAMI is an 'active messaging interface' in that data communications messages in the PAMI are active messages, 'active' in the sense that such messages implement callback functions to advise of message dispatch and instruction completion and so on, thereby reducing the quantity of acknowledgment traffic, and the like, burdening the data communication resources of the PAMI.

Each data communications endpoint of a PAMI is implemented as a combination of a client, a context, and a task. A 'client' as the term is used in PAMI operations is a collection of data communications resources dedicated to the exclusive use of an application-level data processing entity, an application or an application messaging module such as an MPI library. A 'context' as the term is used in PAMI operations is composed of a subset of a client's collection of data processing resources, context functions, and a work queue of data transfer instructions to be performed by use of the subset through the context functions operated by an assigned thread of execution. In at least some embodiments, the context's subset of a client's data processing resources is dedicated to the exclusive use of the context. A 'task' as the term is used in PAMI operations refers to a canonical entity, an integer or objection oriented programming object, that represents in a PAMI a process of execution of the parallel application. That is, a task is typically implemented as an identifier of a particular instance of an application executing on a compute node, a compute core on a compute node, or a thread of execution on a multi-threading compute core on a compute node. In the example of FIG. 1, the compute nodes (102), as well as PAMI endpoints on the compute nodes, are coupled for data communications through the PAMI (218) and through data communications resources such as collective network (106) and point-to-point network (108).

The example parallel computer (100) of FIG. 1 is improved to optimize collective communications according to embodiments of the present invention. Each compute node (102) in

the example of FIG. 1 is configured to execute a plurality of processes. Such a process may be a process in PAMI (a PAMI endpoint, for example), a process representing an instance of an application, or other type of process. Each compute node also includes one or more processors and memory, such as

random access memory (RAM). In typical operation, an operating system within a compute node assigns portions of address space to each processor of the compute node, and, to the extent that the processors include multiple compute cores, treats each compute core as a separate processor with its own assignment of a portion of core memory or RAM for a separate heap, stack, memory variable storage, and so on. The default architecture for such apportionment of memory space is that each processor or compute core operates its assigned portion of memory separately, with no ability to access memory assigned to another processor or compute core. Upon request, however, the operating system grants to one processor or compute core the ability to access a segment of memory that is assigned to another processor or compute core, and such a segment is referred to in this specification as a 'segment of shared memory.'

In the example of FIG. 1, an example processor (199) and RAM (156) of a compute node is illustrated. A collective communication optimizer (198) includes computer program instructions that when executed cause the example processor (199) to carry out the step of determining for each software thread, an affinity of the software thread to a particular hardware thread. A software thread is a thread of execution containing a private stack and local variables in a common address space as other threads. Each affinity indicates an assignment of a software thread to a particular hardware thread. The collective communication optimizer (198) also includes computer program instructions that when executed cause the example processor (199) to carry out the step of generating one or more affinity domains based on the affinities of the software threads. An affinity domain indicates which software threads are assigned to hardware threads of the same hardware domain. A hardware domain may refer to an organizational level within a processor. Examples of a hardware domain include a core domain, a multi-chip module (MCM) domain, and a processor domain. The collective communication optimizer (198) also includes computer program instructions that when executed cause the example processor (199) to carry out the step of generating, for each affinity domain, a topology of the affinity domain based on the affinities of the software threads to the hardware threads. The collective communication optimizer (198) also includes computer program instructions that when executed cause the example processor (199) to carry out the step of performing, based on the generated topologies of the affinity domains, a collective operation on one or more software threads.

The arrangement of compute nodes, networks, and I/O devices making up the example parallel computer illustrated in FIG. 1 are for explanation only, not for limitation of the present invention. Parallel computers capable of optimizing collective communications according to embodiments of the present invention may include additional nodes, networks, devices, and architectures, not shown in FIG. 1, as will occur to those of skill in the art. For ease of explanation, the parallel computer in the example of FIG. 1 is illustrated with only one processor (199) in a compute node. Readers will recognize that compute nodes in parallel computers that optimize collective communications according to embodiments of the present invention can include any number of processors as may occur to those of skill in the art; each compute node in IBM's BlueGene/Q supercomputer, for example, includes 16

application processors and a management processor. The parallel computer (100) in the example of FIG. 1 includes sixteen compute nodes (102); parallel computers that optimize collective communications according to some embodiments of the present invention include thousands of compute nodes. In addition to Ethernet and JTAG, networks in such data processing systems may support many data communications protocols including for example TCP (Transmission Control Protocol), IP (Internet Protocol), and others as will occur to those of skill in the art. Various embodiments of the present invention may be implemented on a variety of hardware platforms in addition to those illustrated in FIG. 1.

Optimizing collective communications within a parallel computer according to embodiments of the present invention is generally implemented on a parallel computer that includes a plurality of compute nodes. In fact, such computers may include thousands of such compute nodes, with a compute node typically executing at least one instance of a parallel application. Each compute node is in turn itself a computer composed of one or more computer processors, its own computer memory, and its own input/output ('I/O') adapters. For further explanation, therefore, FIG. 2 sets forth a block diagram of an example compute node (152) for use in a parallel computer that optimizes collective communications according to embodiments of the present invention. The compute node (152) of FIG. 2 includes one or more computer processors (164) as well as random access memory ('RAM') (156). Each processor (164) can support multiple hardware compute cores (165), and each such core can in turn support multiple threads of execution, hardware threads of execution as well as software threads. Each processor (164) is connected to RAM (156) through a high-speed front side bus (161), bus adapter (194), and a high-speed memory bus (154)—and through bus adapter (194) and an extension bus (168) to other components of the compute node. Stored in RAM (156) is an application program (158) comprising a plurality of software threads (197) having computer program instructions that carry out parallel, user-level data processing using parallel algorithms.

Also stored RAM (156) is an application messaging module (216), a library of computer program instructions that carry out application-level parallel communications among compute nodes, including point to point operations as well as collective operations. Although the application program can call PAMI routines directly, the application program (158) often executes point-to-point data communications operations by calling software routines in the application messaging module (216), which in turn is improved according to embodiments of the present invention to use PAMI functions to implement such communications. An application messaging module can be developed from scratch to use a PAMI according to embodiments of the present invention, using a traditional programming language such as the C programming language or C++, for example, and using traditional programming methods to write parallel communications routines that send and receive data among PAMI endpoints and compute nodes through data communications networks or shared-memory transfers. In this approach, the application messaging module (216) exposes a traditional interface, such as MPI, to the application program (158) so that the application program can gain the benefits of a PAMI with no need to recode the application. As an alternative to coding from scratch, therefore, existing prior art application messaging modules may be improved to use the PAMI, existing modules that already implement a traditional interface. Examples of prior-art application messaging modules that can be improved to optimize collective communications in a parallel computer according to embodiments of the present inven-

tion include such parallel communications libraries as the traditional ‘Message Passing Interface’ (‘MPI’) library, the ‘Parallel Virtual Machine’ (‘PVM’) library, MPICH, and the like.

Also represented in RAM in the example of FIG. 2 is a PAMI (218). Readers will recognize, however, that the representation of the PAMI in RAM is a convention for ease of explanation rather than a limitation of the present invention, because the PAMI and its components, endpoints, clients, contexts, and so on, have particular associations with and inclusions of hardware data communications resources. In fact, the PAMI can be implemented partly as software or firmware and hardware—or even, at least in some embodiments, entirely in hardware.

Also represented in RAM (156) in the example of FIG. 2 is a collective communication optimizer (199). A collective communication optimizer (198) includes computer program instructions that when executed cause at least one of the processors (164) to carry out the step of determining for each software thread, an affinity of the software thread to a particular hardware thread. Each affinity indicates an assignment of a software thread to a particular hardware thread. The collective communication optimizer (198) also includes computer program instructions that when executed cause at least one of the processors (164) to carry out the step of generating one or more affinity domains based on the affinities of the software threads. An affinity domain indicates which software threads are assigned to hardware threads of the same hardware domain. A hardware domain may refer to an organizational level within a processor. Examples of a hardware domain include a core domain, a multi-chip module (MCM) domain, and a processor domain. The collective communication optimizer (198) also includes computer program instructions that when executed cause at least one of the processors (164) to carry out the step of generating, for each affinity domain, a topology of the affinity domain based on the affinities of the software threads to the hardware threads. The collective communication optimizer (198) also includes computer program instructions that when executed cause at least one of the processors (164) to carry out the step of performing, based on the generated topologies of the affinity domains, a collective operation on one or more software threads.

In the example of FIG. 2, each processor or compute core has uniform access to the RAM (156) on the compute node, so that accessing a segment of shared memory is equally fast regardless where the shared segment is located in physical memory. In some embodiments, however, modules of physical memory are dedicated to particular processors, so that a processor may access local memory quickly and remote memory more slowly, a configuration referred to as a Non-Uniform Memory Access or ‘NUMA.’ In such embodiments, a segment of shared memory can be configured locally for one endpoint and remotely for another endpoint—or remotely from both endpoints of a communication. From the perspective of an origin endpoint transmitting data through a segment of shared memory that is configured remotely with respect to the origin endpoint, transmitting data through the segment of shared memory will appear slower than if the segment of shared memory were configured locally with respect to the origin endpoint—or if the segment were local to both the origin endpoint and the target endpoint. This is the effect of the architecture represented by the compute node (152) in the example of FIG. 2 with all processors and all compute cores coupled through the same bus to the RAM—that all accesses to segments of memory shared among processes or processors on the compute node are local—and therefore very fast.

In some embodiments, the application agent is configured to execute on one core (165) of the processors (164), while no other application executes on the same core. In this way, a portion of the shared memory segment (227) may be designated to that core and the application agent’s use.

Also stored in RAM (156) in the example compute node of FIG. 2 is an operating system (162), a module of computer program instructions and routines for an application program’s access to other resources of the compute node. It is possible, in some embodiments at least, for an application program, an application messaging module, and a PAMI in a compute node of a parallel computer to run threads of execution with no user login and no security issues because each such thread is entitled to complete access to all resources of the node. The quantity and complexity of duties to be performed by an operating system on a compute node in a parallel computer therefore can be somewhat smaller and less complex than those of an operating system on a serial computer with many threads running simultaneously with various level of authorization for access to resources. In addition, there is no video I/O on the compute node (152) of FIG. 2, another factor that decreases the demands on the operating system. The operating system may therefore be quite lightweight by comparison with operating systems of general purpose computers, a pared down or ‘lightweight’ version as it were, or an operating system developed specifically for operations on a particular parallel computer. Operating systems that may be improved or simplified for use in a compute node according to embodiments of the present invention include UNIX™, Linux™, Microsoft XP™, AIX™, IBM’s i5/OS™, and others as will occur to those of skill in the art.

The example compute node (152) of FIG. 2 includes several communications adapters (172, 176, 180, 188) for implementing data communications with other nodes of a parallel computer. Such data communications may be carried out serially through RS-232 connections, through external buses such as USB, through data communications networks such as IP networks, and in other ways as will occur to those of skill in the art. Communications adapters implement the hardware level of data communications through which one computer sends data communications to another computer, directly or through a network. Examples of communications adapters for use in computers that optimize collective communications according to embodiments of the present invention include modems for wired communications, Ethernet (IEEE 802.3) adapters for wired network communications, and 802.11b adapters for wireless network communications.

The data communications adapters in the example of FIG. 2 include a Gigabit Ethernet adapter (172) that couples example compute node (152) for data communications to a Gigabit Ethernet (174). Gigabit Ethernet is a network transmission standard, defined in the IEEE 802.3 standard, that provides a data rate of 1 billion bits per second (one gigabit). Gigabit Ethernet is a variant of Ethernet that operates over multimode fiber optic cable, single mode fiber optic cable, or unshielded twisted pair.

The data communications adapters in the example of FIG. 2 includes a JTAG Slave circuit (176) that couples example compute node (152) for data communications to a JTAG Master circuit (178). JTAG is the usual name for the IEEE 1149.1 standard entitled Standard Test Access Port and Boundary-Scan Architecture for test access ports used for testing printed circuit boards using boundary scan. JTAG is so widely adapted that, at this time, boundary scan is more or less synonymous with JTAG. JTAG is used not only for printed circuit boards, but also for conducting boundary scans of integrated circuits, and is also used as a mechanism for



debugging embedded systems. The example compute node of FIG. 2 may be all three of these: It typically includes one or more integrated circuits installed on a printed circuit board and may be implemented as an embedded system having its own processor, its own memory, and its own I/O capability. JTAG boundary scans through JTAG Slave (176) may efficiently configure processor registers and memory in compute node (152) for use in optimizing collective communications according to embodiments of the present invention.

The data communications adapters in the example of FIG. 2 includes a Point To Point Adapter (180) that couples example compute node (152) for data communications to a data communications network (108) that is optimal for point to point message passing operations such as, for example, a network configured as a three-dimensional torus or mesh. Point To Point Adapter (180) provides data communications in six directions on three communications axes, x, y, and z, through six bidirectional links: +x (181), -x (182), +y (183), -y (184), +z (185), and -z (186). For ease of explanation, the Point To Point Adapter (180) of FIG. 2 as illustrated is configured for data communications in three dimensions, x, y, and z, but readers will recognize that Point To Point Adapters optimized for point-to-point operations in a parallel computer that optimizes collective communications according to embodiments of the present invention may in fact be implemented so as to support communications in two dimensions, four dimensions, five dimensions, and so on.

The data communications adapters in the example of FIG. 2 includes a Collective Operations Adapter (188) that couples example compute node (152) for data communications to a network (106) that is optimal for collective message passing operations such as, for example, a network configured as a binary tree. Collective Operations Adapter (188) provides data communications through three bidirectional links: two to children nodes (190) and one to a parent node (192).

The example compute node (152) includes a number of arithmetic logic units ('ALUs'). ALUs (166) are components of processors (164), and a separate ALU (170) is dedicated to the exclusive use of collective operations adapter (188) for use in performing the arithmetic and logical functions of reduction operations. Computer program instructions of a reduction routine in an application messaging module (216) or a PAMI (218) may latch an instruction for an arithmetic or logical function into instruction register (169). When the arithmetic or logical function of a reduction operation is a 'sum' or a 'logical OR,' for example, collective operations adapter (188) may execute the arithmetic or logical operation by use of an ALU (166) in a processor (164) or, typically much faster, by use of the dedicated ALU (170).

The example compute node (152) of FIG. 2 includes a direct memory access ('DMA') controller (225), a module of automated computing machinery that implements, through communications with other DMA engines on other compute nodes, or on a same compute node, direct memory access to and from memory on its own compute node as well as memory on other compute nodes. Direct memory access is a way of reading and writing to and from memory of compute nodes with reduced operational burden on computer processors (164); a CPU initiates a DMA transfer, but the CPU does not execute the DMA transfer. A DMA transfer essentially copies a block of memory from one compute node to another, or between RAM segments of applications on the same compute node, from an origin to a target for a PUT operation, from a target to an origin for a GET operation.

For further explanation, FIG. 3A illustrates an example of a Point To Point Adapter (180) useful in parallel computers that optimize collective communications according to

embodiments of the present invention. Point To Point Adapter (180) is designed for use in a data communications network optimized for point to point operations, a network that organizes compute nodes in a three-dimensional torus or mesh. Point To Point Adapter (180) in the example of FIG. 3A provides data communication along an x-axis through four unidirectional data communications links, to and from the next node in the -x direction (182) and to and from the next node in the +x direction (181). Point To Point Adapter (180) also provides data communication along a y-axis through four unidirectional data communications links, to and from the next node in the -y direction (184) and to and from the next node in the +y direction (183). Point To Point Adapter (180) in also provides data communication along a z-axis through four unidirectional data communications links, to and from the next node in the -z direction (186) and to and from the next node in the +z direction (185). For ease of explanation, the Point To Point Adapter (180) of FIG. 3A as illustrated is configured for data communications in only three dimensions, x, y, and z, but readers will recognize that Point To Point Adapters optimized for point-to-point operations in a parallel computer that optimizes collective communications according to embodiments of the present invention may in fact be implemented so as to support communications in two dimensions, four dimensions, five dimensions, and so on. Several supercomputers now use five dimensional mesh or torus networks, including, for example, IBM's Blue Gene Q™.

For further explanation, FIG. 3B illustrates an example of a Collective Operations Adapter (188) useful in a parallel computer that optimizes collective communications according to embodiments of the present invention. Collective Operations Adapter (188) is designed for use in a network optimized for collective operations, a network that organizes compute nodes of a parallel computer in a binary tree. Collective Operations Adapter (188) in the example of FIG. 3B provides data communication to and from two children nodes through four unidirectional data communications links (190). Collective Operations Adapter (188) also provides data communication to and from a parent node through two unidirectional data communications links (192).

For further explanation, FIG. 4 sets forth a line drawing illustrating an example data communications network (108) optimized for point-to-point operations useful in parallel computers that optimize collective communications according to embodiments of the present invention. In the example of FIG. 4, dots represent compute nodes (102) of a parallel computer, and the dotted lines between the dots represent data communications links (103) between compute nodes. The data communications links are implemented with point-to-point data communications adapters similar to the one illustrated for example in FIG. 3A, with data communications links on three axis, x, y, and z, and to and fro in six directions +x (181), -x (182), +y (183), -y (184), +z (185), and -z (186). The links and compute nodes are organized by this data communications network optimized for point-to-point operations into a three dimensional mesh (105). The mesh (105) has wrap-around links on each axis that connect the outermost compute nodes in the mesh (105) on opposite sides of the mesh (105). These wrap-around links form a torus (107). Each compute node in the torus has a location in the torus that is uniquely specified by a set of x, y, z coordinates. Readers will note that the wrap-around links in the y and z directions have been omitted for clarity, but are configured in a similar manner to the wrap-around link illustrated in the x direction. For clarity of explanation, the data communications network of FIG. 4 is illustrated with only 27 compute nodes, but

readers will recognize that a data communications network optimized for point-to-point operations in a parallel computer that optimizes collective communications according to embodiments of the present invention may contain only a few compute nodes or may contain thousands of compute nodes. For ease of explanation, the data communications network of FIG. 4 is illustrated with only three dimensions: x, y, and z, but readers will recognize that a data communications network optimized for point-to-point operations may in fact be implemented in two dimensions, four dimensions, five dimensions, and so on. As mentioned, several supercomputers now use five dimensional mesh or torus networks, including IBM's Blue Gene Q™.

For further explanation, FIG. 5 illustrates an example data communications network (106) optimized for collective operations by organizing compute nodes in a tree. The example data communications network of FIG. 5 includes data communications links connected to the compute nodes so as to organize the compute nodes as a tree. In the example of FIG. 5, dots represent compute nodes (102) of a parallel computer, and the dotted lines (103) between the dots represent data communications links between compute nodes. The data communications links are implemented with collective operations data communications adapters similar to the one illustrated for example in FIG. 3B, with each node typically providing data communications to and from two children nodes and data communications to and from a parent node, with some exceptions. Nodes in a binary tree may be characterized as a root node (202), branch nodes (204), and leaf nodes (206). The root node (202) has two children but no parent. The leaf nodes (206) each has a parent, but leaf nodes have no children. The branch nodes (204) each has both a parent and two children. The links and compute nodes are thereby organized by this data communications network optimized for collective operations into a binary tree (106). For clarity of explanation, the data communications network of FIG. 5 is illustrated with only 31 compute nodes, but readers will recognize that a data communications network optimized for collective operations for use in parallel computers that optimize collective communications according to embodiments of the present invention may contain only a few compute nodes or hundreds or thousands of compute nodes.

In the example of FIG. 5, each node in the tree is assigned a unit identifier referred to as a 'rank' (196). The rank actually identifies an instance of a parallel application that is executing on a compute node. That is, the rank is an application-level identifier. Using the rank to identify a node assumes that only one such instance of an application is executing on each node. A compute node can, however, support multiple processors, each of which can support multiple processing cores—so that more than one process or instance of an application can easily be present under execution on any given compute node—or in all the compute nodes, for that matter. To the extent that more than one instance of an application executes on a single compute node, the rank identifies the instance of the application as such rather than the compute node. A rank uniquely identifies an application's location in the tree network for use in both point-to-point and collective operations in the tree network. The ranks in this example are assigned as integers beginning with '0' assigned to the root instance or root node (202), '1' assigned to the first node in the second layer of the tree, '2' assigned to the second node in the second layer of the tree, '3' assigned to the first node in the third layer of the tree, '4' assigned to the second node in the third layer of the tree, and so on. For ease of illustration, only the ranks of the first three layers of the tree are shown here, but all compute nodes, or rather all application instances, in the tree network are

assigned a unique rank. Such rank values can also be assigned as identifiers of application instances as organized in a mesh or torus network.

For further explanation, FIG. 6 sets forth a block diagram of an example protocol stack useful in parallel computers that optimize collective communications according to embodiments of the present invention. The example protocol stack of FIG. 6 includes a hardware layer (214), a system messaging layer (212), an application messaging layer (210), and an application layer (208). For ease of explanation, the protocol layers in the example stack of FIG. 6 are shown connecting an origin compute node (222) and a target compute node (224), although it is worthwhile to point out that in embodiments that effect DMA data transfers, the origin compute node and the target compute node can be the same compute node. The granularity of connection through the system messaging layer (212), which is implemented with a PAMI (218), is finer than merely compute node to compute node—because, again, communications among endpoints often is communications among endpoints on the same compute node. For further explanation, recall that the PAMI (218) connects endpoints, connections specified by combinations of clients, contexts, and tasks, each such combination being specific to a thread of execution on a compute node, with each compute node capable of supporting many threads and therefore many endpoints. Every endpoint typically can function as both an origin endpoint or a target endpoint for data transfers through a PAMI, and both the origin endpoint and its target endpoint can be located on the same compute node. So an origin compute node (222) and its target compute node (224) can in fact, and often will, be the same compute node.

The application layer (208) provides communications among instances of a parallel application (158) running on the compute nodes (222, 224) by invoking functions in an application messaging module (216) installed on each compute node. Communications among instances of the application through messages passed between the instances of the application. Applications may communicate messages invoking function of an application programming interface ('API') exposed by the application messaging module (216). In this approach, the application messaging module (216) exposes a traditional interface, such as an API of an MPI library, to the application program (158) so that the application program can gain the benefits of a PAMI, reduced network traffic, callback functions, and so on, with no need to recode the application. Alternatively, if the parallel application is programmed to use PAMI functions, the application can call the PAMI functions directly, without going through the application messaging module.

The example protocol stack of FIG. 6 includes a system messaging layer (212) implemented here as a PAMI (218). The PAMI provides system-level data communications functions that support messaging in the application layer (602) and the application messaging layer (610). Such system-level functions are typically invoked through an API exposed to the application messaging modules (216) in the application messaging layer (210). Although developers can in fact access a PAMI API directly by coding an application to do so, a PAMI's system-level functions in the system messaging layer (212) in many embodiments are isolated from the application layer (208) by the application messaging layer (210), making the application layer somewhat independent of system specific details. With an application messaging module presenting a standard MPI API to an application, for example, with the application messaging module retooled to use the PAMI to carry out the low-level messaging functions, the application gains the benefits of a PAMI with no need to incur the

15

expense of reprogramming the application to call the PAMI directly. Because, however, some applications will in fact be reprogrammed to call the PAMI directly, all entities in the protocol stack above the PAMI are viewed by PAMI as applications. When PAMI functions are invoked by entities above the PAMI in the stack, the PAMI makes no distinction whether the caller is in the application layer or the application messaging layer, no distinction whether the caller is an application as such or an MPI library function invoked by an application. As far as the PAMI is concerned, any caller of a PAMI function is an application.

The protocol stack of FIG. 6 includes a hardware layer (634) that defines the physical implementation and the electrical implementation of aspects of the hardware on the compute nodes such as the bus, network cabling, connector types, physical data rates, data transmission encoding and many other factors for communications between the compute nodes (222) on the physical network medium. In parallel computers that optimize collective communications according to embodiments of the present invention, the hardware layer includes DMA controllers and network links, including routers, packet switches, and the like.

For further explanation, FIG. 7 sets forth a functional block diagram of an example PAMI (218) for use in parallel computers that optimize collective communications according to embodiments of the present invention. The PAMI (218) provides an active messaging layer that supports both point to point communications in a mesh or torus as well as collective operations, gathers, reductions, barriers, and the like in tree networks, for example. The PAMI is a multithreaded parallel communications engine designed to provide low level message passing functions, many of which are one-sided, and abstract such functions for higher level messaging middleware, referred to in this specification as 'application messaging modules' in an application messaging layer. In the example of FIG. 7, the application messaging layer is represented by a generic MPI module (258), appropriate for ease of explanation because some form of MPI is a de facto standard for such messaging middleware. Compute nodes and communications endpoints of a parallel computer (102 on FIG. 1) are coupled for data communications through such a PAMI and through data communications resources (294, 296, 314) that include DMA controllers, network adapters, and data communications networks through which controllers and adapters deliver data communications. The PAMI (218) provides data communications among data communications endpoints, where each endpoint is specified by data communications parameters for a thread of execution on a compute node, including specifications of a client, a context, and a task.

The PAMI (218) in this example includes PAMI clients (302, 304), tasks (286, 298), contexts (190, 292, 310, 312), and endpoints (288, 300). A PAMI client is a collection of data communications resources (294, 295, 314) dedicated to the exclusive use of an application-level data processing entity, an application or an application messaging module such as an MPI library. Data communications resources assigned in collections to PAMI clients are explained in more detail below with reference to FIGS. 8A and 8B. PAMI clients (203, 304 on FIG. 7) enable higher level middleware, application messaging modules, MPI libraries, and the like, to be developed independently so that each can be used concurrently by an application. Although the application messaging layer in FIG. 7 is represented for example by a single generic MPI module (258), in fact, a PAMI, operating multiple clients, can support multiple message passing libraries or application messaging modules simultaneously, a fact that is explained in more

16

detail with reference to FIG. 9. FIG. 9 sets forth a functional block diagram of an example PAMI (218) useful in parallel computers that optimize collective communications according to embodiments of the present invention in which the example PAMI operates, on behalf of an application (158), with multiple application messaging modules (502-510) simultaneously. The application (158) can have multiple messages in transit simultaneously through each of the application messaging modules (502-510). Each context (512-520) carries out, through post and advance functions, data communications for the application on data communications resources in the exclusive possession, in each client, of that context. Each context carries out data communications operations independently and in parallel with other contexts in the same or other clients. In the example FIG. 9, each client (532-540) includes a collection of data communications resources (522-530) dedicated to the exclusive use of an application-level data processing entity, one of the application messaging modules (502-510):

IBM MPI Library (502) operates through context (512) data communications resources (522) dedicated to the use of PAMI client (532),

MPICH Library (504) operates through context (514) data communications resources (524) dedicated to the use of PAMI client (534),

Unified Parallel C ('UPC') Library (506) operates through context (516) data communications resources (526) dedicated to the use of PAMI client (536),

Partitioned Global Access Space ('PGAS') Runtime Library (508) operates through context (518) data communications resources (528) dedicated to the use of PAMI client (538), and

Aggregate Remote Memory Copy Interface ('ARMCI') Library (510) operates through context (520) data communications resources (530) dedicated to the use of PAMI client (540).

Again referring to the example of FIG. 7: The PAMI (218) includes tasks, listed in task lists (286, 298) and identified (250) to the application (158). A 'task' as the term is used in PAMI operations is a platform-defined integer datatype that identifies a canonical application process, an instance of a parallel application (158). Very carefully in this specification, the term 'task' is always used to refer only to this PAMI structure, not the traditional use of the computer term 'task' to refer to a process or thread of execution. In this specification, the term 'process' refers to a canonical data processing process, a container for threads in a multithreading environment. In particular in the example of FIG. 7, the application (158) is implemented as a canonical process with multiple threads (251-254) assigned various duties by a leading thread (251) which itself executes an instance of a parallel application program. Each instance of a parallel application is assigned a task; each task so assigned can be an integer value, for example, in a C environment, or a separate task object in a C++ or Java environment. The tasks are components of communications endpoints, but are not themselves communications endpoints; tasks are not addressed directly for data communications in PAMI. This gives a finer grained control than was available in prior message passing art. Each client has its own list (286, 298) of tasks for which its contexts provide services; this allows each process to potentially reside simultaneously in two or more different communications domains as will be the case in certain advanced computers using, for example, one type of processor and network in one domain and a completely different processor type and network in another domain, all in the same computer.

The PAMI (218) includes contexts (290, 292, 310, 312). A ‘context’ as the term is used in PAMI operations is composed of a subset of a client’s collection of data processing resources, context functions, and a work queue of data transfer instructions to be performed by use of the subset through the context functions operated by an assigned thread of execution. That is, a context represents a partition of the local data communications resources assigned to a PAMI client. Every context within a client has equivalent functionality and semantics. Context functions implement contexts as threading points that applications use to optimize concurrent communications. Communications initiated by a local process, an instance of a parallel application, uses a context object to identify the specific threading point that will be used to issue a particular communication independent of communications occurring in other contexts. In the example of FIG. 7, where the application (158) and the application messaging module (258) are both implemented as canonical processes with multiple threads of execution, each has assigned or mapped particular threads (253, 254, 262, 264) to advance (268, 270, 276, 278) work on the contexts (290, 292, 310, 312), including execution of local callbacks (272, 280). In particular, the local event callback functions (272, 280) associated with any particular communication are invoked by the thread advancing the context that was used to initiate the communication operation in the first place. Like PAMI tasks, contexts are not used to directly address a communication destination or target, as they are a local resource.

Context functions, explained here with regard to references (472-482) on FIG. 9, include functions to create (472) and destroy (474) contexts, functions to lock (476) and unlock (478) access to a context, and functions to post (480) and advance (480) work in a context. For ease of explanation, the context functions (472-482) are illustrated in only one expanded context (512); readers will understand, however, that all PAMI contexts have similar context functions. The create (472) and destroy (474) functions are, in an object-oriented sense, constructors and destructors. In the example embodiments described in this specifications, post (480) and advance (482) functions on a context are critical sections, not thread safe. Applications using such non-reentrant functions must somehow ensure that critical sections are protected from re-entrant use. Applications can use mutual exclusion locks to protect critical sections. The lock (476) and unlock (478) functions in the example of FIG. 9 provide and operate such a mutual exclusion lock to protect the critical sections in the post (480) and advance (482) functions. If only a single thread posts or advances work on a context, then that thread need never lock that context. To the extent that progress is driven independently on a context by a single thread of execution, then no mutual exclusion locking of the context itself is required—provided that no other thread ever attempts to call a function on such a context. If more than one thread will post or advance work on a context, each such thread must secure a lock before calling a post or an advance function on that context. This is one reason why it is probably a preferred architecture, given sufficient resources, to assign one thread to operate each context. Progress can be driven with advance (482) functions concurrently among multiple contexts by using multiple threads, as desired by an application—shown in the example of FIG. 7 by threads (253, 254, 262, 264) which advance work concurrently, independently and in parallel, on contexts (290, 292, 310, 312).

Posts and advances (480, 482 on FIG. 9) are functions called on a context, either in a C-type function with a context ID as a parameter, or in object oriented practice where the calling entity possesses a reference to a context or a context

object as such and the posts and advances are member methods of a context object. Again referring to FIG. 7: Application-level entities, application programs (158) and application messaging modules (258), post (266, 274) data communications instructions, including SENDs, RECEIVES, PUTs, GETs, and so on, to the work queues (282, 284, 306, 308) in contexts and then call advance functions (268, 270, 276, 278) on the contexts to progress specific data processing and data communications that carry out the instructions. The data processing and data communications effected by the advance functions include specific messages, request to send (‘RTS’) messages, acknowledgments, call-back execution, transfers of transfer data or payload data, and so on. Advance functions therefore operate generally by checking a work queue for any new instructions that need to be initiated and checking data communications resources for any incoming message traffic that needs to be administered as well as increases in storage space available for outgoing message traffic, with callbacks and the like. Advance functions also carry out or trigger transfers of transfer data or payload data.

In at least some embodiments, a context’s subset of a client’s data processing resources is dedicated to the exclusive use of the context. In the example of FIG. 7, context (290) has a subset (294) of a client’s (302) data processing resources dedicated to the exclusive use of the context (290), and context (292) has a subset (296) of a client’s (302) data processing resources dedicated to the exclusive use of the context (292). Advance functions (268, 270) called on contexts (290, 292) therefore never need to secure a lock on a data communications resource before progressing work on a context—because each context (290, 292) has exclusive use of dedicated data communications resources. Usage of data communications resources in this example PAMI (218), however, is not thread-safe. When data communications resources are shared among contexts, mutual exclusion locks are needed. In contrast to the exclusive usage of resources by contexts (290, 292), contexts (310, 312) share access to their client’s data communications resource (314) and therefore do not have data communications resources dedicated to exclusive use of a single context. Contexts (310, 312) therefore always must secure a mutual exclusion lock on a data communications resource before using the resource to send or receive administrative messages or transfer data.

For further explanation, here is an example pseudocode Hello World program for an application using a PAMI:

---

```
int main(int argc, char ** argv)
{
    PAMI_client_t client;
    PAMI_context_t context;
    PAMI_result_t status = PAMI_ERROR;
    const char *name = "PAMI";
    status = PAMI_Client_initialize(name, &client);
    size_t _n = 1;
    status = PAMI_Context_create(client, NULL, 0, &context, _n);
    PAMI_configuration_t configuration;
    configuration.name = PAMI_TASK_ID;
    status = PAMI_Configuration_query(client, &configuration);
    size_t task_id = configuration.value.intval;
    configuration.name = PAMI_NUM_TASKS;
    status = PAMI_Configuration_query(client, &configuration);
    size_t num_tasks = configuration.value.intval;
    fprintf(stderr, "Hello process %d of %d\n", task_id, num_tasks);
    status = PAMI_Context_destroy(context);
    status = PAMI_Client_finalize(client);
    return 0;
}
```

---

This short program is termed ‘pseudocode’ because it is an explanation in the form of computer code, not a working model, not an actual program for execution. In this pseudocode example, an application initializes a client and a context for an application named “PAMI.” PAMI\_Client\_initialize and PAMI\_Context\_create are initialization functions (316) exposed to applications as part of a PAMI’s API. These functions, in dependence upon the application name “PAMI,” pull from a PAMI configuration (318) the information needed to establish a client and a context for the application. The application uses this segment:

---

```
PAMI_configuration_t configuration;
configuration.name = PAMI_TASK_ID;
status = PAMI_Configuration_query(client, &configuration);
size_t task_id = configuration.value.intval;
```

---

to retrieve its task ID and this segment:

---

```
configuration.name = PAMI_NUM_TASKS;
status = PAMI_Configuration_query(client, &configuration);
size_t num_tasks = configuration.value.intval;
```

---

to retrieve the number of tasks presently configured to carry out parallel communications and process data communications event in the PAMI. The applications prints “Hello process task\_id of num\_tasks,” where task\_id is the task ID of the subject instance of a parallel application, and num\_tasks is the number of instances of the application executing in parallel on compute nodes. Finally, the application destroys the context and terminates the client.

For further explanation of data communications resources assigned in collections to PAMI clients, FIG. 8A sets forth a block diagram of example data communications resources (220) useful in parallel computers that optimize collective communications according to embodiments of the present invention. The data communications resources of FIG. 8A include a gigabit Ethernet adapter (238), an Infiniband adapter (240), a Fibre Channel adapter (242), a PCI Express adapter (246), a collective operations network configured as a tree (106), shared memory (227), DMA controllers (225, 226), and a network (108) configured as a point-to-point torus or mesh like the network described above with reference to FIG. 4. A PAMI is configured with clients, each of which is in turn configured with certain collections of such data communications resources—so that, for example, the PAMI client (302) in the PAMI (218) in the example of FIG. 7 can have dedicated to its use a collection of data communications resources composed of six segments (227) of shared memory, six Gigabit Ethernet adapters (238), and six Infiniband adapters (240). And the PAMI client (304) can have dedicated to its use six Fibre Channel adapters (242), a DMA controller (225), a torus network (108), and five segments (227) of shared memory. And so on.

The DMA controllers (225, 226) in the example of FIG. 8A each is configured with DMA control logic in the form of a DMA engine (228, 229), an injection FIFO buffer (230), and a receive FIFO buffer (232). The DMA engines (228, 229) can be implemented as hardware components, logic networks of a DMA controller, in firmware, as software operating an embedded controller, as various combinations of software, firmware, or hardware, and so on. Each DMA engine (228, 229) operates on behalf of endpoints to send and receive DMA transfer data through the network (108). The DMA engines (228, 229) operate the injection buffers (230, 232) by

processing first-in-first-out descriptors (234, 236) in the buffers, hence the designation ‘injection FIFO’ and ‘receive FIFO.’

For further explanation, here is an example use case, a description of the overall operation of an example PUT DMA transfer using the DMA controllers (225, 226) and network (108) in the example of FIG. 8A: An originating application (158), which is typically one instance of a parallel application running on a compute node, places a quantity of transfer data (494) at a location in its RAM (155). The application (158) then calls a post function (480) on a context (512) of an origin endpoint (352), posting a PUT instruction (390) into a work queue (282) of the context (512); the PUT instruction (390) specifies a target endpoint (354) to which the transfer data is to be sent as well as source and destination memory locations. The application then calls an advance function (482) on the context (512). The advance function (482) finds the new PUT instruction in its work queue (282) and inserts a data descriptor (234) into the injection FIFO of the origin DMA controller (225); the data descriptor includes the source and destination memory locations and the specification of the target endpoint. The origin DMA engine (225) then transfers the data descriptor (234) as well as the transfer data (494) through the network (108) to the DMA controller (226) on the target side of the transaction. The target DMA engine (229), upon receiving the data descriptor and the transfer data, places the transfer data (494) into the RAM (156) of the target application at the location specified in the data descriptor and inserts into the target DMA controller’s receive FIFO (232) a data descriptor (236) that specifies the target endpoint and the location of the transfer data (494) in RAM (156). The target application (159) or application instance calls an advance function (483) on a context (513) of the target endpoint (354). The advance function (483) checks the communications resources assigned to its context (513) for incoming messages, including checking the receive FIFO (232) of the target DMA controller (226) for data descriptors that specify the target endpoint (354). The advance function (483) finds the data descriptor for the PUT transfer and advises the target application (159) that its transfer data has arrived. A GET-type DMA transfer works in a similar manner, with some differences, including, of course, the fact that transfer data flows in the opposite direction. Similarly, typical SEND transfers also operate similarly, some with rendezvous protocols, some with eager protocols, with data transmitted in packets over the a network through non-DMA network adapters or through DMA controllers.

The example of FIG. 8A includes two DMA controllers (225, 226). DMA transfers between endpoints on separate compute nodes use two DMA controllers, one on each compute node. Compute nodes can be implemented with multiple DMA controllers so that many or even all DMA transfers even among endpoints on a same compute node can be carried out using two DMA engines. In some embodiments at least, however, a compute node, like the example compute node (152) of FIG. 2, has only one DMA engine, so that that DMA engine can be used to conduct both sides of transfers between endpoints on that compute node. For further explanation of this fact, FIG. 8B sets forth a functional block diagram of an example DMA controller (225) operatively coupled to a network (108)—in an architecture where this DMA controller (225) is the only DMA controller on a compute node—and an origin endpoint (352) and its target endpoint (354) are both located on the same compute node (152). In the example of FIG. 8B, a single DMA engine (228) operates with two threads of execution (502, 504) on behalf of endpoints (352, 354) on a same compute node to send and receive DMA

transfer data through a segment (227) of shared memory. A transmit thread (502) injects transfer data into the network (108) as specified in data descriptors (234) in an injection FIFO buffer (230), and a receive thread (502) receives transfer data from the network (108) as specified in data descriptors (236) in a receive FIFO buffer (232).

The overall operation of an example PUT DMA transfer with the DMA controllers (225) and the network (108) in the example of FIG. 8B is: An originating application (158), that is actually one of multiple instances (158, 159) of a parallel application running on a compute node (152) in separate threads of execution, places a quantity of transfer data (494) at a location in its RAM (155). The application (158) then calls a post function (480) on a context (512) of an origin endpoint (352), posting a PUT instruction (390) into a work queue (282) of the context (512); the PUT instruction specifies a target endpoint (354) to which the transfer data is to be sent as well as source and destination memory locations. The application (158) then calls an advance function (482) on the context (512). The advance function (482) finds the new PUT instruction (390) in its work queue (282) and inserts a data descriptor (234) into the injection FIFO of the DMA controller (225); the data descriptor includes the source and destination memory locations and the specification of the target endpoint. The DMA engine (225) then transfers by its transmit and receive threads (502, 504) through the network (108) the data descriptor (234) as well as the transfer data (494). The DMA engine (228), upon receiving by its receive thread (504) the data descriptor and the transfer data, places the transfer data (494) into the RAM (156) of the target application and inserts into the DMA controller's receive FIFO (232) a data descriptor (236) that specifies the target endpoint and the location of the transfer data (494) in RAM (156). The target application (159) calls an advance function (483) on a context (513) of the target endpoint (354). The advance function (483) checks the communications resources assigned to its context for incoming messages, including checking the receive FIFO (232) of the DMA controller (225) for data descriptors that specify the target endpoint (354). The advance function (483) finds the data descriptor for the PUT transfer and advises the target application (159) that its transfer data has arrived. Again, a GET-type DMA transfer works in a similar manner, with some differences, including, of course, the fact that transfer data flows in the opposite direction. And typical SEND transfers also operate similarly, some with rendezvous protocols, some with eager protocols, with data transmitted in packets over the a network through non-DMA network adapters or through DMA controllers.

By use of an architecture like that illustrated and described with reference to FIG. 8B, a parallel application or an application messaging module that is already programmed to use DMA transfers can gain the benefit of the speed of DMA data transfers among endpoints on the same compute node with no need to reprogram the applications or the application messaging modules to use the network in other modes. In this way, an application or an application messaging module, already programmed for DMA, can use the same DMA calls through a same API for DMA regardless whether subject endpoints are on the same compute node or on separate compute nodes.

For further explanation, FIG. 10 sets forth a functional block diagram of example endpoints useful in parallel computers that optimize collective communications according to embodiments of the present invention. In the example of FIG. 10, a PAMI (218) is implemented with instances on two separate compute nodes (152, 153) that include four endpoints (338, 340, 342, 344). These endpoints are opaque

objects used to address an origin or destination in a process and are constructed from a (client, task, context) tuple. Non-DMA SEND and RECEIVE instructions as well as DMA instructions such as PUT and GET address a destination by use of an endpoint object or endpoint identifier.

Each endpoint (338, 340, 342, 344) in the example of FIG. 10 is composed of a client (302, 303, 304, 305), a task (332, 333, 334, 335), and a context (290, 292, 310, 312). Using a client a component in the specification of an endpoint disambiguates the task and context identifiers, as these identifiers may be the same for multiple clients. A task is used as a component in the specification of an endpoint to construct an endpoint to address a process accessible through a context. A context in the specification of an endpoint identifies, refers to, or represents the specific context associated with a destination or target task—because the context identifies a specific threading point on a task. A context offset identifies which threading point is to process a particular communications operation. Endpoints enable “crosstalk” which is the act of issuing communication on a local context with a particular context offset that is directed to a destination endpoint with no correspondence to a source context or source context offset.

For efficient utilization of storage in an environment where multiple tasks of a client reside on the same physical compute node, an application may choose to write an endpoint table (288, 300 on FIG. 7) in a segment of shared memory (227, 346, 348). It is the responsibility of the application to allocate such segments of shared memory and coordinate the initialization and access of any data structures shared between processes. This includes any endpoint objects which are created by one process or instance of an application and read by another process.

Endpoints (342, 344) on compute node (153) serve respectively two application instances (157, 159). The tasks (334, 336) in endpoints (342, 344) are different. The task (334) in endpoint (342) is identified by the task ID (249) of application (157), and the task (336) in endpoint (344) is identified by the task ID (257) of application (159). The clients (304, 305) in endpoints (342, 344) are different, separate clients. Client (304) in endpoint (342) associates data communications resources (e.g., 294, 296, 314 on FIG. 7) dedicated exclusively to the use of application (157), while client (305) in endpoint (344) associates data communications resources dedicated exclusively to the use of application (159). Contexts (310, 312) in endpoints (342, 344) are different, separate contexts. Context (310) in endpoint (342) operates on behalf of application (157) a subset of the data communications resources of client (304), and context (312) in endpoint (344) operates on behalf of application (159) a subset of the data communications resources of client (305).

Contrasted with the PAMIs (218) on compute node (153), the PAMI (218) on compute node (152) serves only one instance of a parallel application (158) with two endpoints (338, 340). The tasks (332, 333) in endpoints (338, 340) are the same, because they both represent a same instance of a same application (158); both tasks (332, 333) therefore are identified, either with a same variable value, references to a same object, or the like, by the task ID (250) of application (158). The clients (302, 303) in endpoints (338, 340) are optionally either different, separate clients or the same client. If they are different, each associates a separate collection of data communications resources. If they are the same, then each client (302, 303) in the PAMI (218) on compute node (152) associates a same set of data communications resources and is identified with a same value, object reference, or the like. Contexts (290, 292) in endpoints (338, 340) are different, separate contexts. Context (290) in endpoint (338) oper-

ates on behalf of application (158) a subset of the data communications resources of client (302) regardless whether clients (302, 303) are the same client or different clients, and context (292) in endpoint (340) operates on behalf of application (158) a subset of the data communications resources of client (303) regardless whether clients (302, 303) are the same client or different clients. Thus the tasks (332, 333) are the same; the clients (302, 303) can be the same; and the endpoints (338, 340) are distinguished at least by different contexts (290, 292), each of which operates on behalf of one of the threads (251-254) of application (158), identified typically by a context offset or a threading point.

Endpoints (338, 340) being as they are on the same compute node (152) can effect DMA data transfers between endpoints (338, 340) through DMA controller (225) and a segment of shared local memory (227). In the absence of such shared memory (227), endpoints (338, 340) can effect DMA data transfers through the DMA controller (225) and the network (108), even though both endpoints (338, 340) are on the same compute node (152). DMA transfers between endpoint (340) on compute node (152) and endpoint (344) on another compute node (153) go through DMA controllers (225, 226) and either a network (108) or a segment of shared remote memory (346). DMA transfers between endpoint (338) on compute node (152) and endpoint (342) on another compute node (153) also go through DMA controllers (225, 226) and either a network (108) or a segment of shared remote memory (346). The segment of shared remote memory (346) is a component of a Non-Uniform Memory Access ('NUMA') architecture, a segment in a memory module installed anywhere in the architecture of a parallel computer except on a local compute node. The segment of shared remote memory (346) is 'remote' in the sense that it is not installed on a local compute node. A local compute node is 'local' to the endpoints located on that particular compute node. The segment of shared remote memory (346), therefore, is 'remote' with respect to endpoints (338, 340) on compute node (158) if it is in a memory module on compute node (153) or anywhere else in the same parallel computer except on compute node (158).

Endpoints (342, 344) being as they are on the same compute node (153) can effect DMA data transfers between endpoints (342, 344) through DMA controller (226) and a segment of shared local memory (348). In the absence of such shared memory (348), endpoints (342, 344) can effect DMA data transfers through the DMA controller (226) and the network (108), even though both endpoints (342, 344) are on the same compute node (153). DMA transfers between endpoint (344) on compute node (153) and endpoint (340) on another compute node (152) go through DMA controllers (226, 225) and either a network (108) or a segment of shared remote memory (346). DMA transfers between endpoint (342) on compute node (153) and endpoint (338) on another compute node (158) go through DMA controllers (226, 225) and either a network (108) or a segment of shared remote memory (346). Again, the segment of shared remote memory (346) is 'remote' with respect to endpoints (342, 344) on compute node (153) if it is in a memory module on compute node (158) or anywhere else in the same parallel computer except on compute node (153).

FIG. 11 sets forth a flow chart illustrating a functional block diagram of an example processor (199) that includes a collective communication optimizer (198) useful for optimization of collective communications in a parallel computer according to embodiments of the present invention.

In the example of FIG. 11, the processor (199) includes four multi-chip modules (MCMs) (1110-1113). A multi-chip

module (MCM) is a specialized electronic package where multiple integrated circuits (ICs), semiconductor dies or other discrete components are packaged onto a unifying substrate, facilitating their use as a single component (as though a larger IC). The MCM itself will often be referred to as a "chip" in designs, thus illustrating its integrated nature. Although only four MCMs are illustrated in the example of FIG. 11, readers of skill in the art will realize that a plurality of any number of MCMs may be included in the processor (199). Each MCM in the example of FIG. 11 includes four cores. A core or central processing unit (CPU) is hardware within a computer system which carries out the instructions of a computer program by performing the basic arithmetical, logical, and input/output operations of the system. For example, first MCM (1110) includes cores (1120-1123), second MCM (1111) includes cores (1124-1127), third MCM (1112) includes cores (1128-1131), and fourth MCM (1113) includes cores (1132-1135). Although only four cores are illustrated in each MCM, readers of skill in the art will realize that a plurality of any number of cores may be included in each MCM. In addition, each core includes a plurality of hardware threads (1160-1175). A hardware thread is an execution block that shares resources of a single core: the computing units, the CPU caches and the translation lookaside buffer (TLB). Although only four hardware threads are illustrated in each core, readers of skill in the art will realize that a plurality of any number of hardware threads may be included in each core.

For further explanation, FIG. 12 sets forth a flow chart illustrating an example method of optimizing collective communications within a parallel computer according to embodiments of the present invention. For ease of explanation, the method of FIG. 12 is explained with reference to elements, components, and processes of FIGS. 1-11. The method of FIG. 12 is carried out in a parallel computer similar to the example parallel computer (100) of FIG. 1. Such a parallel computer (100) includes a plurality of compute nodes. Each compute node is configured to execute a plurality of processes and parallel application (156).

The method of FIG. 12 includes determining (1202) for each software thread (197), an affinity (1220) of the software thread (197) to a particular hardware thread (1160-1175). Each affinity indicates an assignment of a software thread to a particular hardware thread. Determining (1202) an affinity (1220) of the software thread (197) to a particular hardware thread (1160-1175) may be carried out by retrieving location information of software threads and using the location information to determine which software threads are pinned or assigned to a particular hardware thread. For example, a first software thread may be assigned to one of the hardware threads (1160) in a first core (1120) and a second software thread may be assigned to one of the hardware threads (1161) in a second core (1121).

The method of FIG. 12 also includes generating (1204) one or more affinity domains (1222) based on the affinities (1220) of the software threads (197). An affinity domain indicates which software threads are assigned to hardware threads of the same hardware domain. A hardware domain may refer to an organizational level within a processor. Examples of hardware domains include a core domain, a multi-chip module (MCM) domain, and a processor domain. Generating (1204) one or more affinity domains (1222) may be carried out by determining which software threads are assigned to hardware threads of a particular hardware domain and creating a list containing indications of all of the software threads determined to be assigned to hardware threads of that hardware domain. For example, an affinity domain that indicates 'core' affinity would include indications of all the software threads



25

assigned to hardware threads within that particular core. As another example, an affinity domain that indicates MCM affinity would include indications of all the software threads assigned to hardware threads within that particular MCM.

The method of FIG. 12 also includes generating (1206), for each affinity domain (1222), a topology (1224) of the affinity domain based on the affinities (1220) of the software threads (197) to the hardware threads (1160-1175). Generating (1206) a topology (1224) of an affinity domain may be carried out by determining communication connections for software threads associated with a particular affinity domain. For example, an affinity domain corresponding to MCM (1110) of FIG. 1 may include software threads assigned to the plurality of hardware threads (1160) of the core (1120) and software threads assigned to the plurality of hardware threads (1161) of the core (1121). In this example, a topology generated based on the affinities (1120) may create one topology for communications between software threads assigned to the core (1120), another topology for communications between software threads assigned to the core (1121), and a third topology for communications between a software thread assigned to the core (1120) and a software thread assigned to the core (1121).

The method of FIG. 12 also includes performing (1208) a collective operation (1226) on one or more software threads (197) based on the generated topologies (1224) of the affinity domains (1222). Collective operations are implemented with data communications among the compute nodes of an operational group. Collective operations are those functions that involve all the compute nodes of an operational group. A collective operation is an operation, a message-passing computer program instruction that is executed simultaneously, that is, at approximately the same time, by all the compute nodes in an operational group of compute nodes. Performing (1208) a collective operation (1226) on one or more software threads (197) based on the generated topologies (1224) of the affinity domains (1222) may be carried out by sending and receiving data between the software threads in accordance with the generated topologies (1224). For example, software threads associated with a core affinity domain may transmit data based on a n-ary tree pattern while software threads from different affinity domains wishing to communicate may use designated software threads (or affinity domain leaders) to communicate between affinity domains.

For further explanation, FIG. 13 sets forth a flow chart illustrating an example method of optimizing collective communications within a parallel computer according to embodiments of the present invention. For ease of explanation, the method of FIG. 13 is explained with reference to elements, components, and processes of FIGS. 1-11. The method of FIG. 13 is carried out in a parallel computer similar to the example parallel computer (100) of FIG. 1. Such a parallel computer (100) includes a plurality of hardware threads for executing software threads of a parallel application.

The method of FIG. 13 is similar to the method of FIG. 12 in that the method of FIG. 13 includes: determining (1202) for each software thread (197), an affinity (1220) of the software thread (197) to a particular hardware thread (1160-1175); generating (1204), based on the affinities (1220) of the software threads (197), one or more affinity domains (1222); generating (1206), for each affinity domain (1222), a topology (1224) of the affinity domain based on the affinities (1220) of the software threads (197) to the hardware threads (1160-1175); and performing (1208), based on the generated topologies (1224) of the affinity domains (1222), a collective operation (1226) on one or more software threads (197).

26

In the method of FIG. 13, however, generating (1204) one or more affinity domains (1222) includes generating (1302), for each core (1120-1127), a core affinity domain (1320) indicating the software threads assigned to the hardware threads within the core. Generating (1302) a core affinity domain (1320) may be carried out by determining which software threads are assigned to hardware threads of a particular core and creating a list contains indications of all of the software threads determined to be assigned to hardware threads of that core.

In the method of FIG. 13, generating (1204) one or more affinity domains (1222) also includes generating (1304), for each MCM (1110-1113), a MCM affinity domain (1322) indicating the software threads assigned to the hardware threads within the MCM. Generating (1304) a MCM affinity domain (1322) may be carried out by determining which software threads are assigned to hardware threads of a particular MCM and creating a list contains indications of all of the software threads determined to be assigned to hardware threads of that MCM.

In the method of FIG. 13, generating (1204) one or more affinity domains (1222) also includes generating (1306), for the processor (199), a processor affinity domain (1324) indicating the software threads assigned to the hardware threads within the processor (199). Generating (1306) a processor affinity domain (1324) may be carried out by determining which software threads are assigned to hardware threads of a particular processor and creating a list contains indications of all of the software threads determined to be assigned to hardware threads of that processor.

For further explanation, FIG. 14 sets forth a flow chart illustrating an example method of optimizing collective communications within a parallel computer according to embodiments of the present invention. For ease of explanation, the method of FIG. 14 is explained with reference to elements, components, and processes of FIGS. 1-11. The method of FIG. 14 is carried out in a parallel computer similar to the example parallel computer (100) of FIG. 1. Such a parallel computer (100) includes a plurality of hardware threads for executing software threads of a parallel application.

The method of FIG. 14 is similar to the method of FIG. 12 in that the method of FIG. 14 includes: determining (1202) for each software thread (197), an affinity (1220) of the software thread (197) to a particular hardware thread (1160-1175); generating (1204), based on the affinities (1220) of the software threads (197), one or more affinity domains (1222); generating (1206), for each affinity domain (1222), a topology (1224) of the affinity domain based on the affinities (1220) of the software threads (197) to the hardware threads (1160-1175); and performing (1208), based on the generated topologies (1224) of the affinity domains (1222), a collective operation (1226) on one or more software threads (197).

In the method of FIG. 14, however, determining (1202) for each software thread (197), an affinity (1220) of the software thread (197) to a particular hardware thread (1160-1175) includes querying (1402) an operating system (162) for software thread locality information (1420) corresponding to the software thread. Querying (1402) an operating system (162) for software thread locality information (1420) corresponding to the software thread may be carried out by sending a request to the operating system (162) to retrieve software thread locality information and receiving the software thread locality information.

In the method of FIG. 14, generating (1206), for each affinity domain (1222), a topology (1224) of the affinity domain based on the affinities (1220) of the software threads (197) to the hardware threads (1160-1175) includes generat-



ing (1404), for each affinity domain, an n-ary tree (1422) representing a communication organization among the software threads associated with the affinity domain. An n-ary tree is a rooted tree in which each node has no more than n children. Generating (1404), for each affinity domain, an n-ary tree (1422) representing a communication organization among the software threads associated with the affinity domain may be carried out by creating a data structure indicating which software threads communicate with each other in a particular affinity domain.

In the method of FIG. 14, performing (1208), based on the generated topologies (1224) of the affinity domains (1222), a collective operation (1226) on one or more software threads (197) includes performing (1406), for each affinity domain, in accordance with the generated n-ary tree (1422), a reduction operation (1424) on the software threads associated with the affinity domain. A collective operation requires that all processes on all compute nodes within an operational group call the same collective operation with matching arguments. A 'reduce' operation is an example of a collective operation that executes arithmetic or logical functions on data distributed among the compute nodes of an operational group. Performing (1406) a reduction operation (1424) on the software threads associated with the affinity domain in accordance with the generated n-ary tree (1422) may be carried out by carrying out arithmetic or logical functions on data from software threads within the affinity domain.

For further explanation, FIG. 15 sets forth a flow chart illustrating an example method of optimizing collective communications within a parallel computer according to embodiments of the present invention. For ease of explanation, the method of FIG. 15 is explained with reference to elements, components, and processes of FIGS. 1-11. The method of FIG. 15 is carried out in a parallel computer similar to the example parallel computer (100) of FIG. 1. Such a parallel computer (100) includes a plurality of hardware threads for executing software threads of a parallel application.

The method of FIG. 15 is similar to the method of FIG. 12 in that the method of FIG. 15 includes: determining (1202) for each software thread (197), an affinity (1220) of the software thread (197) to a particular hardware thread (1160-1175); generating (1204), based on the affinities (1220) of the software threads (197), one or more affinity domains (1222); generating (1206), for each affinity domain (1222), a topology (1224) of the affinity domain based on the affinities (1220) of the software threads (197) to the hardware threads (1160-1175); and performing (1208), based on the generated topologies (1224) of the affinity domains (1222), a collective operation (1226) on one or more software threads (197).

In the method of FIG. 15, generating (1206), for each affinity domain (1222), a topology (1224) of the affinity domain based on the affinities (1220) of the software threads (197) to the hardware threads (1160-1175) includes designating (1502), for each affinity domain, an affinity domain leader (1522). Designating (1502) an affinity domain leader (1522) may be carried out by selecting one of the software threads associated with a particular affinity domain to be the 'leader' of that affinity domain. The affinity domain leader may be used to communicate with other affinity domain leaders such that software threads in different affinity domains may use only the domain leaders to communicate with each other.

In the method of FIG. 15, performing (1208), based on the generated topologies (1224) of the affinity domains (1222), a collective operation (1226) on one or more software threads (197) includes performing (1504) a barrier operation (1524) on each of the affinity domain leaders (1522). A collective operation requires that all processes on all compute nodes

within an operational group call the same collective operation with matching arguments. A 'barrier operation' is an example of a collective operation that is issued by each thread participating in the barrier operation and only can complete when all threads have entered and left the collective barrier operation. Performing (1504) a barrier operation (1524) on each of the affinity domain leaders (1522) may be carried out by each affinity domain leader issuing a barrier operation and waiting until each of the other affinity domain leaders has entered and left the barrier operation.

In the method of FIG. 15, performing (1208), based on the generated topologies (1224) of the affinity domains (1222), a collective operation (1226) on one or more software threads (197) includes broadcasting (1506), for each affinity domain leader (1522), from the affinity domain leader to the software threads within the affinity domain corresponding to the affinity domain leader, results (1526) of the barrier operation (1524). A collective operation requires that all processes on all compute nodes within an operational group call the same collective operation with matching arguments. A 'broadcast' is an example of a collective operation for moving data among compute nodes of an operational group. Broadcasting (1506) from the affinity domain leader to the software threads within the affinity domain corresponding to the affinity domain leader, results (1526) of the barrier operation (1524) may be carried out by transmitting between affinity domain leaders, the results of the barrier operations.

As will be appreciated by one skilled in the art, aspects of the present invention may be embodied as a system, method or computer program product. Accordingly, aspects of the present invention may take the form of an entirely hardware embodiment, an entirely software embodiment (including firmware, resident software, micro-code, etc.) or an embodiment combining software and hardware aspects that may all generally be referred to herein as a "circuit," "module" or "system." Furthermore, aspects of the present invention may take the form of a computer program product embodied in one or more computer readable medium(s) having computer readable program code embodied thereon.

Any combination of one or more computer readable medium(s) may be utilized. The computer readable medium may be a computer readable signal medium or a computer readable storage medium. A computer readable storage medium may be, for example, but not limited to, an electronic, magnetic, optical, electromagnetic, infrared, or semiconductor system, apparatus, or device, or any suitable combination of the foregoing. More specific examples (a non-exhaustive list) of the computer readable storage medium would include the following: an electrical connection having one or more wires, a portable computer diskette, a hard disk, a random access memory (RAM), a read-only memory (ROM), an erasable programmable read-only memory (EPROM or Flash memory), an optical fiber, a portable compact disc read-only memory (CD-ROM), an optical storage device, a magnetic storage device, or any suitable combination of the foregoing. In the context of this document, a computer readable storage medium may be any tangible medium that can contain, or store a program for use by or in connection with an instruction execution system, apparatus, or device.

A computer readable signal medium may include a propagated data signal with computer readable program code embodied therein, for example, in baseband or as part of a carrier wave. Such a propagated signal may take any of a variety of forms, including, but not limited to, electro-magnetic, optical, or any suitable combination thereof. A computer readable signal medium may be any computer readable medium that is not a computer readable storage medium and

that can communicate, propagate, or transport a program for use by or in connection with an instruction execution system, apparatus, or device.

Program code embodied on a computer readable medium may be transmitted using any appropriate medium, including but not limited to wireless, wireline, optical fiber cable, RF, etc., or any suitable combination of the foregoing.

Computer program code for carrying out operations for aspects of the present invention may be written in any combination of one or more programming languages, including an object oriented programming language such as Java, Smalltalk, C++ or the like and conventional procedural programming languages, such as the "C" programming language or similar programming languages. The program code may execute entirely on the user's computer, partly on the user's computer, as a stand-alone software package, partly on the user's computer and partly on a remote computer or entirely on the remote computer or server. In the latter scenario, the remote computer may be connected to the user's computer through any type of network, including a local area network (LAN) or a wide area network (WAN), or the connection may be made to an external computer (for example, through the Internet using an Internet Service Provider).

Aspects of the present invention are described above with reference to flowchart illustrations and/or block diagrams of methods, apparatus (systems) and computer program products according to embodiments of the invention. It will be understood that each block of the flowchart illustrations and/or block diagrams, and combinations of blocks in the flowchart illustrations and/or block diagrams, can be implemented by computer program instructions. These computer program instructions may be provided to a processor of a general purpose computer, special purpose computer, or other programmable data processing apparatus to produce a machine, such that the instructions, which execute via the processor of the computer or other programmable data processing apparatus, create means for implementing the functions/acts specified in the flowchart and/or block diagram block or blocks.

These computer program instructions may also be stored in a computer readable medium that can direct a computer, other programmable data processing apparatus, or other devices to function in a particular manner, such that the instructions stored in the computer readable medium produce an article of manufacture including instructions which implement the function/act specified in the flowchart and/or block diagram block or blocks.

The computer program instructions may also be loaded onto a computer, other programmable data processing apparatus, or other devices to cause a series of operational steps to be performed on the computer, other programmable apparatus or other devices to produce a computer implemented process such that the instructions which execute on the computer or other programmable apparatus provide processes for implementing the functions/acts specified in the flowchart and/or block diagram block or blocks.

The flowchart and block diagrams in the Figures illustrate the architecture, functionality, and operation of possible implementations of systems, methods and computer program products according to various embodiments of the present invention. In this regard, each block in the flowchart or block diagrams may represent a module, segment, or portion of code, which comprises one or more executable instructions for implementing the specified logical function(s). It should also be noted that, in some alternative implementations, the functions noted in the block may occur out of the order noted in the figures. For example, two blocks shown in succession

may, in fact, be executed substantially concurrently, or the blocks may sometimes be executed in the reverse order, depending upon the functionality involved. It will also be noted that each block of the block diagrams and/or flowchart illustration, and combinations of blocks in the block diagrams and/or flowchart illustration, can be implemented by special purpose hardware-based systems that perform the specified functions or acts, or combinations of special purpose hardware and computer instructions.

It will be understood from the foregoing description that modifications and changes may be made in various embodiments of the present invention without departing from its true spirit. The descriptions in this specification are for purposes of illustration only and are not to be construed in a limiting sense. The scope of the present invention is limited only by the language of the following claims.

What is claimed is:

1. A method of optimizing collective communications within a parallel computer, the parallel computer comprising a plurality of hardware threads for executing software threads of a parallel application, the method comprising:

determining for each software thread, by a processor of the parallel computer, an affinity of the software thread to a particular hardware thread, each affinity indicating an assignment of a software thread to a particular hardware thread, wherein the processor further comprises one or more multi-chip modules (MCM) each MCM comprising a plurality of cores;

generating, based on the affinities of the software threads, one or more affinity domains, wherein an affinity domain indicates which software threads are assigned to hardware threads of a same hardware domain, including: generating, for each core, a core affinity domain indicating the software threads assigned to the hardware threads within the core,

generating, for each MCM, a MCM affinity domain indicating the software threads assigned to the hardware threads within the MCM, and

generating, for the processor, a processor affinity domain indicating the software threads assigned to the hardware threads within the processor;

generating, for each affinity domain, a topology of the affinity domain based on the affinities of the software threads to the hardware threads, including generating, for each affinity domain, an n-ary tree representing a communication organization among the software threads associated with the affinity domain; and

performing, based on the generated topologies of the affinity domains, a collective operation on one or more software threads, wherein performing the collective operation on one or more software threads based on the generated topologies of the affinity domains includes performing, for each affinity domain, in accordance with the generated n-ary tree, a reduction operation on the software threads associated with the affinity domain.

2. The method of claim 1 wherein generating, for each affinity domain, the topology of the affinity domain based on the affinities of the software threads to the hardware threads includes designating, for each affinity domain, an affinity domain leader.

3. The method of claim 2 wherein performing the collective operation on one or more software threads based on the generated topologies of the affinity domains includes:

performing a barrier operation on each of the affinity domain leaders; and

broadcasting, for each affinity domain leader, from the affinity domain leader to the software threads within the

31

affinity domain corresponding to the affinity domain leader, results of the barrier operation.

4. The method of claim 1, wherein determining for each software thread, by the processor of the parallel computer, an affinity of the software thread to a particular hardware thread includes querying an operating system for software thread locality information corresponding to the software thread.

5. An apparatus for optimizing collective communications within a parallel computer, the parallel computer comprising a plurality of hardware threads for executing software threads of a parallel application, the apparatus comprising a computer processor, a computer memory operatively coupled to the computer processor, the computer memory having disposed within it computer program instructions that, when executed, cause the apparatus to carry out the steps of:

determining for each software thread, by a processor of the parallel computer, an affinity of the software thread to a particular hardware thread, each affinity indicating an assignment of a software thread to a particular hardware thread, wherein the processor further comprises one or more multi-chip modules (MCM), each MCM comprising a plurality of cores;

generating, based on the affinities of the software threads, one or more affinity domains, wherein an affinity domain indicates which software threads are assigned to hardware threads of a same hardware domain, including: generating, for each core, a core affinity domain indicating the software threads assigned to the hardware threads within the core,

generating, for each MCM, a MCM affinity domain indicating the software threads assigned to the hardware threads within the MCM, and

generating, for the processor, a processor affinity domain indicating the software threads assigned to the hardware threads within the processor;

generating, for each affinity domain, a topology of the affinity domain based on the affinities of the software threads to the hardware threads, including generating, for each affinity domain, an n-ary tree representing a communication organization among the software threads associated with the affinity domain; and

performing, based on the generated topologies of the affinity domains, a collective operation on one or more software threads, wherein performing the collective operation on one or more software threads based on the generated topologies of the affinity domains includes performing, for each affinity domain, in accordance with the generated n-ary tree, a reduction operation on the software threads associated with the affinity domain.

6. The apparatus of claim 5 wherein generating, for each affinity domain, the topology of the affinity domain based on the affinities of the software threads to the hardware threads includes designating, for each affinity domain, an affinity domain leader.

7. The apparatus of claim 6 wherein performing the collective operation on one or more software threads based on the generated topologies of the affinity domains includes:

performing a barrier operation on each of the affinity domain leaders; and

broadcasting, for each affinity domain leader, from the affinity domain leader to the software threads within the affinity domain corresponding to the affinity domain leader, results of the barrier operation.

8. The apparatus of claim 5, wherein determining for each software thread, by the processor of the parallel computer, an affinity of the software thread to a particular hardware thread

32

includes querying an operating system for software thread locality information corresponding to the software thread.

9. A computer program product for optimizing collective communications within a parallel computer, the parallel computer comprising a plurality of hardware threads for executing software threads of a parallel application, the computer program product disposed upon a computer readable medium, wherein the computer readable medium is not a signal, the computer program product comprising computer program instructions that, when executed, cause a computer to carry out the steps of:

determining for each software thread, by a processor of the parallel computer, an affinity of the software thread to a particular hardware thread, each affinity indicating an assignment of a software thread to a particular hardware thread, wherein the processor further comprises one or more multi-chip modules (MCM) each MCM comprising a plurality of cores;

generating, based on the affinities of the software threads, one or more affinity domains, wherein an affinity domain indicates which software threads are assigned to hardware threads of a same hardware domain, including: generating, for each core, a core affinity domain indicating the software threads assigned to the hardware threads within the core,

generating, for each MCM, a MCM affinity domain indicating the software threads assigned to the hardware threads within the MCM, and

generating, for the processor, a processor affinity domain indicating the software threads assigned to the hardware threads within the processor;

generating, for each affinity domain, a topology of the affinity domain based on the affinities of the software threads to the hardware threads, including generating, for each affinity domain, an n-ary tree representing a communication organization among the software threads associated with the affinity domain; and

performing, based on the generated topologies of the affinity domains, a collective operation on one or more software threads, wherein performing the collective operation on one or more software threads based on the generated topologies of the affinity domains includes performing, for each affinity domain, in accordance with the generated n-ary tree, a reduction operation on the software threads associated with the affinity domain.

10. The computer program product of claim 9, wherein generating, for each affinity domain, the topology of the affinity domain based on the affinities of the software threads to the hardware threads includes designating, for each affinity domain, an affinity domain leader.

11. The computer program product of claim 10, wherein performing the collective operation on one or more software threads based on the generated topologies of the affinity domains includes:

performing a barrier operation on each of the affinity domain leaders; and

broadcasting, for each affinity domain leader, from the affinity domain leader to the software threads within the affinity domain corresponding to the affinity domain leader, results of the barrier operation.

12. The computer program product of claim 9, wherein determining for each software thread, by the processor of the parallel computer, an affinity of the software thread to a particular hardware thread includes querying an operating system for software thread locality information corresponding to the software thread.

13. The computer program product of claim 9 wherein the computer readable medium comprises a computer readable storage medium.

\* \* \* \* \*